

Clemson University

**TigerPrints**

---

All Dissertations

Dissertations

---

August 2020

## Essays on Tactical and Operational Problems in Healthcare

Mahsa Kiani

*Clemson University*, [mkiani@g.clemson.edu](mailto:mkiani@g.clemson.edu)

Follow this and additional works at: [https://tigerprints.clemson.edu/all\\_dissertations](https://tigerprints.clemson.edu/all_dissertations)

---

### Recommended Citation

Kiani, Mahsa, "Essays on Tactical and Operational Problems in Healthcare" (2020). *All Dissertations*. 2688.

[https://tigerprints.clemson.edu/all\\_dissertations/2688](https://tigerprints.clemson.edu/all_dissertations/2688)

This Dissertation is brought to you for free and open access by the Dissertations at TigerPrints. It has been accepted for inclusion in All Dissertations by an authorized administrator of TigerPrints. For more information, please contact [kokeefe@clemson.edu](mailto:kokeefe@clemson.edu).

# ESSAYS ON TACTICAL AND OPERATIONAL PROBLEMS IN HEALTHCARE

---

A Dissertation  
Presented to  
the Graduate School of  
Clemson University

---

In Partial Fulfillment  
of the Requirements for the Degree  
Doctor of Philosophy  
Industrial Engineering

---

by  
Mahsa Kiani  
August 2020

---

Accepted by:  
Dr. Tugce Isik, Committee Chair  
Dr. Caglar Caglayan  
Dr. Burak Eksioglu  
Dr. Sandra Eksioglu  
Dr. Amin Khademi

# Abstract

One of the essential challenges in healthcare operations management is to efficiently utilize the expensive resources needed in the healthcare system, while maintaining or increasing the quality of care. Optimization methods can be used to increase the supply of healthcare services, to minimize the cost of the system, and maximize the quality of care by minimizing patients' waiting times, minimizing travel needs, maximizing health outcomes and maximizing access to services. In this dissertation, we study some of the important tactical and operational problems in healthcare, and propose plans to efficiently improve the current healthcare systems by applying optimization methods.

In chapter 1 of this dissertation, we develop a novel scheduling model called “postponement model” to reduce the indirect waiting time of higher priority outpatients in a diagnostic clinic. In diagnostic clinics, the arrivals mostly arise from three sources: inpatients, emergency patients, and outpatients. Emergency patients are seen as soon as they arrive and inpatients receive appointments within 24 hours. However, outpatient appointments are scheduled within a longer time horizon based on appointment availability. Currently, most diagnostic clinics save a proportion of their capacity for inpatients and emergency patients, and offer the earliest remaining appointments to the outpatients on a first-come-first-serve basis. This capacity allocation and scheduling mechanism may lead to unused inpatient capacity. Furthermore, there is no prioritization in scheduling of outpatients whose medical needs may be at different urgency levels. We model the appointment scheduling problem as a two-stage stochastic integer program. In the first stage we compute the proportion of capacity that is allocated to emergency patients and inpatients. In the second stage the decisions regarding scheduling of outpatients are taken. Outpatient appointments are not necessarily scheduled immediately upon patients' arrivals and may be postponed to observe more requests. This postponement strategy enables the scheduler to observe more of the demand and schedule outpatient

appointments considering the patient priorities. We solve the problem using Sample Average Approximation (SAA) and a decomposition based branch and bound algorithm. The results show that using the postponement acceptance patients with higher priority receive sooner appointments compared to the no-postponement scheduling policy used in current practice. Meanwhile, the utilization of the system is increased.

In chapter 2, we study a dynamic model for Tuberculosis (TB) screening of healthcare personnel. Healthcare employees take TB diagnostic tests regularly as part of efforts to prevent TB outbreaks in hospitals. A simple strategy that is mostly used in countries with low rate of TB infections is annual screening of all employees. There are currently two TB diagnostic tests on the market: skin test and blood test. The blood test is more expensive than the skin test, however it is more accurate. In this study, we propose an alternative testing scheme where testing frequency and test type for different groups of employees is dependent on their infection risk and the cost of time lost due to testing. We develop a discrete time infinite horizon Markov Decision Process (MDP) model which determines the optimal time between the tests for different groups of employees. Another outcome of our model is the type of the TB diagnostic test administered for each employee group. Classification of employees into groups is done based on the characteristics that affect the probability of getting infected with TB (e.g., job type and work location) and employee salary levels. The objective of our model is to minimize the total cost of the healthcare facility which depends on the type of the tests administered, employees' lost time, and the number of false-positive or false-negative results in each group tested. Due to the curse of dimensionality, we use Approximate Dynamic Programming (ADP) to estimate the value function. Then, we use column generation to solve the ADP-based linear program associated with the proposed MDP model. The results provide screening policies that determine which test should be allocated to each group of employee in different states of the system. By investigating the results, we also estimate the frequency of the test for each group. Comparison of the screening policies obtained using our model with the current annual screening policy show that the screening costs can be reduced by half while achieving the same the overall infection rate among the healthcare personnel.

In chapter 3, we propose a dynamic model for scheduling of healthcare workers during an infectious disease outbreak, with a specific focus on the ongoing Coronavirus Disease 2019 (COVID-19) pandemic. Healthcare workers play an important role during a pandemic to control the infection spread in the general population. On the other hand, they are at high risk of getting infected because

of being in direct contact with patients. Thus, taking operational measures to limit the exposure of healthcare workers to infectious patients is critical for the safety of the healthcare workers and their patients. Emerging literature indicates that creating teams of healthcare workers and scheduling or isolating these teams in coordination might be beneficial during the COVID-19 pandemic. In this study, we build a MDP model to determine the optimal policy for scheduling such worker teams. The objective of the model is to maximize the expected total discounted number of working employees while taking the possibility of infection, and thus quarantine, for workers who are scheduled to work into account. The optimal policy specifies which teams of workers should work and which teams should isolate dependent on the system state. This problem is difficult to solve due to the large size of the state space of the MDP. Thus, we use state space reduction techniques to decrease the number of states. Using the data on number of infections in the state of South Carolina, we obtain optimal scheduling policies under different infection probabilities for the general population. We also consider additional scenarios to understand the effect of changing model parameters on the state space reduction results and the approximate optimal policy. The results show that strategic benching of healthcare worker teams can significantly improve the total discounted workable physician days compared to only segregating workers into teams.

# Dedication

I wish to dedicate this thesis to my late father, Mehrdad. He taught me to persevere and prepared me to face the challenges with faith and humility. He was constant source of inspiration to my life. Although, he is not here to give me strength and support, I always feel his presence that used to urge me to strive to achieve my goals in life.

# Acknowledgments

I am grateful to my loving parent Soudabeh and my wonderful brothers Mahyar and Mohammadreza who have always inspired me through love, encouragement and support. I wish to express my gratitude to my advisor Dr. Tugce Isik and my co-advisor Dr. Burak Eksioglu who believed in me and supported me so I can reach this point as a PhD student. I would also like to thank my committee members, Dr. Caglar Caglayan, Dr. Sandra Eksioglu, and Dr. Amin Khademi for providing valuable insights during my studies as a PhD student, for serving as my committee members, and for their constructive comments and suggestions. I am also grateful to the staff in Industrial Engineering department for their unfailing support and assistance. Finally, I wish to express my gratitude to my friends who have supported me in different stages of my life.

# Table of Contents

|  |           |
|--|-----------|
| <b>Title Page</b> . . . . .  | <b>i</b>  |
| <b>Abstract</b> . . . . .  | <b>ii</b> |
| <b>Dedication</b> . . . . .  | <b>v</b>  |
| <b>Acknowledgments</b> . . . . .   | <b>vi</b> |
| <b>List of Tables</b> . . . . .  | <b>ix</b> |
| <b>List of Figures</b> . . . . .   | <b>x</b>  |
| <b>1 Evaluating Appointment Postponement in Scheduling Patients at a Diagnostic Clinic</b> . . . . . | <b>1</b>  |
| 1.1 Introduction . . . . .   | 2         |
| 1.2 Literature Review . . . . .  | 4         |
| 1.3 Problem Definition and Formulation . . . . .   | 6         |
| 1.4 Solution Approach . . . . .  | 10        |
| 1.5 Numerical Study . . . . .  | 14        |
| 1.6 Conclusion . . . . .   | 24        |
| <b>2 Dynamic Tuberculosis Screening for Healthcare Employees</b> . . . . .                           | <b>26</b> |
| 2.1 Introduction . . . . .   | 26        |
| 2.2 Literature Review . . . . .  | 28        |
| 2.3 An MDP Model for the TB Test Scheduling Problem . . . . .  | 30        |
| 2.4 Approximate Dynamic Programming . . . . .  | 36        |
| 2.5 Numerical Study and Results . . . . .  | 39        |
| 2.6 Conclusion . . . . .   | 44        |
| <b>3 Risk Based Staffing for Pandemic Response</b> . . . . .   | <b>46</b> |
| 3.1 Introduction . . . . .   | 47        |
| 3.2 Literature Review . . . . .  | 48        |
| 3.3 An MDP Model for Staff Scheduling . . . . .  | 50        |
| 3.4 Solution Approach . . . . .  | 55        |
| 3.5 Numerical Study and Numerical Results . . . . .  | 57        |
| 3.6 Conclusion . . . . .   | 65        |
| <b>4 Contributions and Future Research</b> . . . . .   | <b>67</b> |
| <b>Appendices</b> . . . . .  | <b>69</b> |
| A . . . . .  | 70        |
| B . . . . .  | 71        |



|                        |    |
|------------------------|----|
| Bibliography . . . . . | 72 |
|------------------------|----|

# List of Tables

|     |  |    |
|-----|--|----|
| 1.1 | Problem Parameters . . . . .   | 7  |
| 1.2 | Problem Variables . . . . .  | 8  |
| 1.3 | Cost improvement and capacity allocation for the base scenario . . . . .             | 17 |
| 1.4 | Summary of all results for all four policies for the base scenario . . . . .         | 18 |
| 1.5 | Percentage of outpatients waiting in the acceptance queue for one vs. two days . . . | 23 |
| 1.6 | Capacity utilization of each patient type for policy 2 . . . . .                     | 24 |
| 2.1 | Optimal actions based on the simulation . . . . .                                    | 44 |
| 2.2 | Comparison of the proposed policies and the current policy . . . . .                 | 44 |
| 3.1 | Estimation of parameters for each scenario . . . . .                                 | 60 |
| 3.2 | Comparison of the original and reduced problems' sizes . . . . .                     | 60 |
| 3.3 | Comparison of the approximate optimal policy versus the benchmark policy . . . . .   | 64 |
| 1   | Parameters values for the postponement model . . . . .                               | 70 |
| 2   | Estimated parameters in no-postponement model . . . . .                              | 70 |
| 3   | Parameters values . . . . .  | 71 |

# List of Figures

|     |  |    |
|-----|--|----|
| 1.1 | Effect of changing $b_{jt^{u_{ta}}}$ on average total cost and capacity allocation . . . . .                       | 20 |
| 1.2 | Effect of changing $b_{jt^{u_{ta}}}$ on the number of days in acceptance queue and indirect waiting time . . . . . | 21 |
| 2.1 | Optimal policy for risk group 1 . . . . .  | 41 |
| 2.2 | Optimal policy for risk group 2 . . . . .  | 42 |
| 2.3 | Optimal policy for risk group 3 . . . . .  | 43 |
| 3.1 | Daily new cases in South Carolina . . . . .  | 58 |
| 3.2 | Forecast for the daily new cases in South Carolina using New York data . . . . .                                   | 58 |
| 3.3 | Smoothed data for the number of new cases . . . . .  | 59 |
| 3.4 | Smoothed data for the number of new cases with specified $p^t$ values . . . . .                                    | 61 |
| 3.5 | Optimal policy for scenario 1 when $w_1^t = 1, w_2^t = 1$ . . . . .  | 62 |
| 3.6 | Optimal policy for scenario 2 when $w_1^t = 1, w_2^t = 1$ . . . . .  | 63 |
| 3.7 | Optimal policy for scenario 3 when $w_1^t = 1, w_2^t = 1$ . . . . .  | 63 |
| 3.8 | Optimal policy for scenario 6 when $w_1^t = 1, w_2^t = 1, n_3^t = 2$ . . . . .                                     | 64 |

# Chapter 1

## Evaluating Appointment Postponement in Scheduling Patients at a Diagnostic Clinic

**Summary:** Diagnostic clinics are among healthcare facilities that suffer from long waiting times which can cause medical issues and lead to increases in patient no-shows. Reducing waiting times without significant capital investments is a challenging task. We tackle this challenge by proposing a new appointment scheduling model for diagnostics clinics that does not require significant investments. The clinic in our study serves outpatients, inpatients, and emergency patients. Emergency patients must be seen on arrival, and inpatients must be given next day appointments. Outpatients, however, can be given later appointments. The proposed model takes advantage of this flexibility by allowing the postponement of the acceptance of appointment requests from outpatients. The appointment scheduling process is modeled as a two-stage stochastic programming problem where a portion of the clinic capacity is allocated to inpatients and emergency patients in the first stage. In the second stage, outpatients are scheduled based on their priority classes. After a detailed analysis of the solutions obtained from the two-stage stochastic model, we develop a simple, non-anticipative policy for patient scheduling. We evaluate the performance of this proposed, easy-to-implement policy in a simulation study which shows significant improvements in outpatient indirect waiting times.

## 1.1 Introduction

In today’s healthcare systems, the increasing demand for appointments combined with a shortage of physicians has led to challenges for healthcare providers to give timely appointments to patients. To achieve good medical outcomes, offering timely appointments is important [39]. [39] classify waiting time of patients into two categories. They define *direct waiting time* as the time the patient waits in the healthcare facility on the day of appointment and *indirect waiting time* as the time between the day the patient requests an appointment and the appointment day. Unfortunately, long indirect waiting times are common in practice. For instance, [47] reported that 84% of patients in Veterans Affairs (VA) hospitals wait more than 14 days to see a physician. In addition to the medical issues that long indirect waiting times cause, they can also lead to increases in patient no-shows [38] which have significant effect on annual revenues [63]. Thus, healthcare managers face the challenge of improving their appointment systems to decrease waiting times and no-shows without incurring major capital costs.

Diagnostic clinics are among the healthcare facilities that generally suffer from long indirect waiting times [61]. One such clinic is the Radiology Department at Prisma Health, our collaborator on this study. The clinic provides service to outpatients, inpatients, and emergency patients. The requests for appointments are handled on a first-come-first-served (FCFS) basis. The emergency patients are the highest priority group, followed by inpatients and then outpatients. The outpatients are further categorized into a number of priority classes based on co-morbidities and chronic conditions. The emergency patients are seen as soon as they arrive if there is capacity or immediately referred to another clinic. The inpatients are either given a next day appointment during regular hours or seen during overtime hours. The clinic prefers to offer appointments to outpatients within a few days. However, under the current system, the average indirect waiting time for outpatients is about one week. [58] provide other examples where indirect waiting times for outpatients are negatively impacted by the arrival of higher priority inpatients and emergency patients. A possible strategy to reduce the indirect waiting times for outpatients is to allocate a part of the overall capacity for emergency patients to dampen their impact on the overall system. Similarly, a portion of the capacity can also be reserved for inpatients. However, this strategy can result in unused capacity. Meanwhile, the limited available capacity may not allow providers to serve some of the more urgent outpatients in an acceptable time period. Thus, finding ways to utilize the unused portion of

the capacity reserved for inpatients and allocating just enough capacity for emergency patients are important problems.

The clinic currently makes all acceptance and referral decisions upon the arrival of appointment requests. This causes some high priority outpatients to be referred to other clinics while some of the capacity reserved for inpatients goes unused. As a solution, we propose postponing the acceptance of outpatient requests. In other words, the decision regarding acceptance or referral of an outpatient is not taken upon arrival of an appointment request but is revisited after the inpatient schedules are realized. This postponement will enable the scheduling of higher priority arrivals sooner and also allow for better utilization of the unused capacity reserved for inpatients. Note that postponement does not allow one to utilize the potential unused capacity allocated for emergency patients, because we do not have the one day buffer which is the case for inpatients. Thus, it is critical to allocate the right amount of capacity for emergency patients.

The majority of the outpatients prefer to get an immediate response from the clinic regarding their appointment request. However, the clinic is willing to keep outpatient appointment requests in an *acceptance queue* for a reasonable amount of time. While some patients may leave for an alternative healthcare facility, the clinic believes that most of the outpatients will be amenable to waiting in the acceptance queue if it means their total indirect waiting time will be shorter. Still, the clinic is not open to keeping the outpatients in the acceptance queue more than 72 hours.

To that end, we develop a two-stage, postponable acceptance appointment model which first allocates the total regular-time capacity among different groups of patients and then schedules appointments. Outpatient appointment requests are either scheduled during regular hours or referred to another clinic. The objective is to minimize the expected total cost over the planning horizon. The remainder of the study is organized into five additional sections. Section 1.2 provides a review of the relevant literature. In Section 1.3 the problem is formally defined and a notation is provided along with a two-stage stochastic programming (TSSP) formulation. Section 1.4 explains how the problem is solved. Specifically, the details of our sample average approximation (SAA) and decomposition-based branch-and-bound (DBB) algorithm are provided. Section 1.5 shows the results of our extensive experiments and sensitivity analysis. Finally, Section 2.6 concludes the study with some managerial insights, highlights some of the limitations, and provides directions for future research.

## 1.2 Literature Review

Our study is related to four streams of literature. In the following paragraphs we provide brief reviews of the related literature on (i) patient scheduling, (ii) acceptance postponement, (iii) solution approaches for TSSP, and (iv) revenue management. We highlight how our study differs from those in the literature and summarize our contributions.

The scheduling of patients with different priority classes and medical resource allocation to these classes has gotten a lot of attention in recent years, as evidenced by the large number of papers in the literature [70, 73, 17, 33, 52, 43]. [4] provide a comprehensive review of recent analytical and numerical studies in the area of outpatient scheduling. Some of these studies consider inpatients and emergency patients in addition to outpatients, where the arrival of inpatients and emergency patients are modeled as random events that interrupt the system [69, 31, 32]. [27] provide a capacity allocation plan to minimize the indirect waiting time of higher priority patients across an integrated network of care services. On the other hand, scheduling of outpatients in the presence of emergency and inpatient arrivals is studied via appointment scheduling, but not capacity planning, in diagnostic clinics by [38, 78], and [19]. [38] discuss scheduling of patients in a diagnostic clinic where a certain number of outpatients are already scheduled. They assume that emergency patients arrive randomly throughout the day, and they have to be seen as soon as they arrive. They specify which patient to schedule next when both outpatients and inpatients are waiting for appointments. [19] study outpatient and inpatient scheduling problems with non-homogeneous mean service times considering punctuality and no-show rates. Reserving a part of capacity for emergency arrivals or inpatients is shown to be beneficial to decrease the waiting time of urgent patients [69]. [86] apply a robust optimization model in a surgery department to decide how much capacity to allocate for elective surgeries and emergency surgeries when the demand is uncertain. The decision regarding acceptance or rejection of patients depends on their priorities and available capacity. [9] develop a finite-horizon Markov decision process to schedule appointments considering choice behavior and no-show rate of patients. Patients provide their preference for a specific physician and time of appointment. The decision to accept or reject walk-in patients is based on already scheduled patients who called-in. The main difference between our study and those mentioned above is the timing of the decisions regarding acceptance, rejection, or referral of outpatient appointment requests.

The concept of acceptance postponement is developed and discussed in some manufacturing

settings but not so much in service settings. For example, [45] present a model for a manufacturing system with postponable acceptance and assignment in make-to-order settings, where postponement is applied to both acceptance and assignment. In their model, acceptance of some orders may be deferred or cancelled to wait for more profitable orders. They show that by applying this model the total profit of the system improves. In a study by [36], some low-priority orders are rejected or the acceptance decision is postponed to reserve inventory for higher priority orders. [16] provide one of the few studies of applying postponement in a service system. They consider a call center routing problem that assigns arrivals right after acceptance or after some waiting period. However, acceptance of calls have to be made at the time of arrival. Moreover, acceptance and assignment decisions are made at the same time if there is an available agent. The two closest studies to ours are by [15] and [70]. [15] consider both open-access and prescheduled appointments in their settings. They compute how much of a physician’s workload should be allocated to prescheduled appointments. However, scheduling of patients occurs upon their arrivals. In contrast, the study by [70] considers the acceptance of some of the requests to be postponed. They consider a dynamic system which schedules multiple priority classes of outpatients with the goal of decreasing indirect waiting times when the daily outpatient capacity is fixed. In their model, once the acceptance decisions are made, the remaining requests are deferred to the next day and may be accepted later. However, they do not keep track of the number of days that the decisions are deferred. We postpone the acceptance and scheduling of outpatients in our setting as well. However, our study considers the following concepts that are not included in [70] study. First, we consider the cost of postponing the acceptance decisions, which depends on the amount of time outpatients wait in the acceptance queue and their priority classes. Second, we consider an abandonment probability which relies on the outpatient’s priority class and the amount of time they have waited in the acceptance queue. Third, we consider how the postponing of outpatients affects the capacity allocation and scheduling of inpatients and emergency patients. To the best of our knowledge, our study is the first one that introduces a capacity allocation and postponement model for patient scheduling.

As discussed in Section 1.3, we formulate our problem as a TSSP and develop SAA and DBB approaches, as detailed in Section 1.4, to solve the problem. A well-known approach to solve TSSP is stochastic Benders decomposition which is also known as the L-shaped method. [87] were the first ones to use the L-shaped method to solve TSSP problems. In their formulation, the first and second stage variables were all continuous. [54] allowed integer first and/or second stage variables



in their setting by incorporating a branch-and-bound procedure. [7] proposed the DBB algorithm by branching on tender variables that are the product of first stage variables with the technology matrix for problems with integer variables in the second stage. In our study, we first replace the original objective function by a SAA function and then apply the DBB algorithm to be able to solve realistic size problems.

While our study does not directly contribute to the revenue management literature, there are similarities. Revenue management is defined as the management of perishable assets [90]. Examples of perishable assets include hotel rooms, rental cars, and airplane seats. Revenue management of these perishable assets includes the process of allocating a fixed capacity to the right customer at the right time at the right price [80]. One of the studies which is close to ours is where they allocate the scarce inventory to stochastic demand for multiple fare classes so as to maximize the total expected revenue [18]. The structure of optimal policy is estimated by solving an approximate dynamic program. Revenue management decisions are made upon arrivals but considering anticipated future requests. In this perspective, our study is different since we consider the possibility of postponing the decisions.

### 1.3 Problem Definition and Formulation

As mentioned in Section 1.1, the diagnostic clinic in our study receives appointment requests from outpatients, inpatients, and emergency patients. Currently, almost all of the outpatient appointment requests are accepted or referred to another clinic as soon as the request arrives. The only exception to this are those requests that are received via fax which constitute a small fraction of all requests. The clinic responds to the faxed requests by the end of the business day. We, on the other hand, develop a TSSP that allows the postponement of all outpatient requests.

Outpatients are categorized into  $J$  priority classes ( $j = 1, \dots, J$ ). Parts of the regular-time capacity are allocated for inpatients and emergency patients. The capacity reserved for inpatients can be used for outpatients only if it is unused after inpatients are scheduled. Emergency patients that arrive throughout the day are either seen upon arrival or immediately referred to another clinic. Inpatients that arrive throughout the day are either given a next day appointment during regular hours upon arrival or seen during overtime hours. Outpatient requests that arrive each day are kept in the acceptance queue. In other words, the acceptance and scheduling decisions of lower priority

outpatients can be postponed while waiting for inpatients, emergency patients, or higher priority outpatients. To facilitate the formulation of our model we define the parameters shown in Table 1.1 and the variables shown in Table 1.2.

| Parameters   |
|--|
| $T$ : length of the planning horizon ( $t = 1, 2, \dots, T$ )  |
| $T^a$ : length of the booking horizon ( $t^a = 1, 2, \dots, T^a$ )   |
| $T^u$ : maximum number of days an outpatient waits in the acceptance queue ( $t^u = 1, 2, \dots, T^u$ )  |
| $K$ : daily regular-time capacity of the clinic  |
| $p_{jtu}$ : proportion of type $j$ outpatients who stay in the acceptance queue one more day after having waited for $(t^u-1)$ days                        |
| $a_{jtu}$ : cost of a type $j$ outpatient leaving the acceptance queue after waiting for $t^u$ days  |
| $b_{jtu t^a}$ : cost of giving an appointment to a type $j$ outpatient $t^a$ days later when the patient has waited for $t^u$ days in the acceptance queue |
| $c_{jtu}^O$ : cost of referring a type $j$ outpatient to another clinic when the patient has waited for $t^u$ days in the acceptance queue                 |
| $c^I$ : cost of seeing an inpatient during overtime hours  |
| $c^E$ : cost of referring an emergency patient to another clinic   |

Table 1.1: Problem Parameters

In our proposed system, a scheduler observes the number of inpatients ( $D_t^I$ ) and outpatients ( $D_{jt}^O$ ) that have arrived during the day and the available capacity in each future day of the booking horizon. If the daily inpatient arrivals exceed the allocated capacity ( $K\alpha^I$ ), they are handled during overtime hours which incurs additional cost ( $c^I$ ). If any of the capacity allocated to inpatients is not used then it can be allocated to an outpatient from the acceptance queue. However, the capacity reserved for emergency patients ( $K\alpha^E$ ) is never used for inpatients or outpatients. If an emergency patient arrives when the allocated capacity is full then they are immediately referred to another clinic. An outpatient who has been in the acceptance queue for  $T^u$  days is referred to another clinic.

Based on analysis of historical data and our conversations with Prisma Health, patient arrivals are independent Poisson processes. Thus, we model  $D_{jt}^O$ ,  $D_t^I$  and  $D_t^E$  as truncated Poisson distributions with rates  $\lambda_j, \lambda^I$  and  $\lambda^E$ , respectively. The evolution of  $Q_{jtt^u}$ , the number of outpatients in the acceptance queue, is captured by the following equations:

$$\begin{aligned}
Q_{jt1} &= p_{j1} D_{jt}^O - \sum_{t^a=1}^{T^a} y_{j t 1 t^a}^O - r_{j t 1}, & \forall j, t, & \quad (1.1a) \\
Q_{j(t+1)(t^u+1)} &= p_{j(t^u+1)} Q_{j t t^u} - \sum_{t^a=1}^{T^a} y_{j(t+1)(t^u+1) t^a}^O - r_{j(t+1)(t^u+1)}, & \forall j, t, t^u, t \neq T, t^u \neq T^u, t(1 \leq t) &
\end{aligned}$$

|                    |   |
|--------------------|---|
| Random Variables   |   |
| $D_{jt}^O$         | : number of type $j$ outpatients that arrive at the clinic during day $t$   |
| $D_t^I$            | : number of inpatients that arrive at the clinic during day $t$   |
| $D_t^E$            | : number of emergency patients that arrive at the clinic during day $t$   |
| Decision Variables |   |
| $\alpha^I$         | : percentage of total capacity $K$ reserved for inpatients  |
| $\alpha^E$         | : percentage of total capacity $K$ reserved for emergency patients  |
| $y_{jtt^u t^a}^O$  | : number of type $j$ outpatients who are given an appointment in day $t$ for $t^a$ days later after waiting for $t^u$ days in the acceptance queue ( $t^u \leq t$ ) |
| $r_{jtt^u}$        | : number of type $j$ outpatients who are referred to another clinic in day $t$ after waiting for $t^u$ days in the acceptance queue ( $t^u \leq t$ )                |
| $Q_{jtt^u}$        | : number of type $j$ outpatients in day $t$ who have been waiting in the acceptance queue for $t^u$ days ( $t^u \leq t$ )   |
| $K_{tt^a}^O$       | : available capacity for outpatients $t^a$ days after day $t$   |

Table 1.2: Problem Variables

where equation (1.1a) is a special case of equation (1.1b) with  $t^u = 1$ . These two equations simply state that the number of outpatients in the next day will be equal to the number of outpatients who are not scheduled or referred yet and remained in the queue for one more day.

We also need to maintain an accurate account of the remaining regular-time capacity. This can be achieved by the following equations where (1.2a) is for the beginning of the planning horizon, (1.2b) for the end of the booking horizon, and equation (1.2c) for other days during the planning and booking horizons. At the beginning of the planning horizon and the end of the booking horizon we have full capacity for outpatients since no body is scheduled in these days yet. In the remaining days, the available capacity on day  $(t+1)$  is available capacity of day  $t$  minus the scheduled appointments for that day.

$$K_{1t^a}^O = K(1 - \alpha^I - \alpha^E), \quad \forall t^a, \quad (1.2a)$$

$$K_{tT^a}^O = K(1 - \alpha^I - \alpha^E), \quad \forall t, \quad (1.2b)$$

$$K_{(t+1)t^a}^O = K_{t(t^a+1)}^O - \sum_{j=1}^J \sum_{t^u=1}^{T^u} y_{jtt^u(t^a+1)}^O, \quad \forall t, t^a, t \neq T, t^a \neq T^a. \quad (1.2c)$$

The postponable acceptance appointment system can now be formulated as the following TSSP, named (2SIP). Since capacity allocations have to be made prior to the realization of patient arrivals,  $\alpha = (\alpha^I, \alpha^E)$  are the first-stage decision variables. On the other hand, the appointments depend on

patient arrivals. Thus,  $\mathbf{y}^O, \mathbf{r}, \mathbf{Q}$ , and  $\mathbf{K}^O$  are the second-stage variables.

$$(2\text{SIP}) \quad C^* = \min_{\boldsymbol{\alpha}} \mathbb{E}_{\boldsymbol{\omega} \in \Omega} [C(\boldsymbol{\alpha}, \boldsymbol{\omega})] \quad (1.3a)$$

$$\text{s.t.} \quad \alpha^I + \alpha^E \leq 1, \quad (1.3b)$$

$$\alpha^I, \alpha^E \geq 0. \quad (1.3c)$$

The model minimizes the expected total cost associated with appointment scheduling. Note that  $\boldsymbol{\omega} = \{(D_{1t}^O, \dots, D_{Jt}^O, D_t^I, D_t^E) \text{ for } t = 1, \dots, T\}$  is a joint scenario for the planning horizon. We assume that there is no cost for capacity allocation. The objective of the second stage is to minimize the cost associated with scheduling patient appointments. As shown in Table 1.1, costs are incurred when outpatients abandon the acceptance queue, outpatients are given late appointments, outpatients are referred to another clinic, inpatients are seen during overtime hours, and emergency patients are referred to another clinic. Thus,  $C(\boldsymbol{\alpha}, \boldsymbol{\omega})$  is the objective function value of the second-stage problem given below:

$$C(\boldsymbol{\alpha}, \boldsymbol{\omega}) = \min_{\mathbf{y}^O, \mathbf{r}, \mathbf{Q}, \mathbf{K}^O} \left\{ \sum_{t=1}^T \left( \sum_{j=1}^J \sum_{t^u=1}^{T^u} \sum_{t^a=1}^{T^a} b_{jtt^ut^a} y_{jtt^ut^a}^O + \sum_{j=1}^J \sum_{t^u=1}^{T^u} c_{jtu}^O r_{jtt^u} + a_{jtu} (1 - p_{jtu}) Q_{jtt^u} \right) \right. \quad (1.4a)$$

$$\left. + c^I (1 - z_t^I) (D_t^I - \alpha^I K) + c^E (1 - z_t^E) (D_t^E - \alpha^E K) \right\} \quad (1.4b)$$

$$\text{s.t.} \quad (1.1a) - (1.2c), \quad (1.4b)$$

$$z_t^I K \geq \alpha^I K - D_t^I, \quad \forall t, \quad (1.4c)$$

$$z_t^I D_t^I \leq \alpha^I K, \quad \forall t, \quad (1.4d)$$

$$z_t^E K \geq \alpha^E K - D_t^E, \quad \forall t, \quad (1.4e)$$

$$z_t^E D_t^E \leq \alpha^E K, \quad \forall t, \quad (1.4f)$$

$$\sum_{j=1}^J \sum_{t^u=1}^{T^u} y_{jtt^u}^O - z_t^I (\alpha^I K - D_t^I) \leq K_{t1}^O, \quad \forall t, \quad (1.4g)$$

$$\sum_{j=1}^J \sum_{t^u=1}^{T^u} y_{jtt^ut^a}^O \leq K_{tt^a}^O, \quad \forall t, t^a = 2, \dots, T^a, \quad (1.4h)$$

$$Q_{jtt^u} = r_{jtt^u}, \quad \forall j, t, t^u \geq T^u, \quad (1.4i)$$

$$z_t^I, z_t^E \in \{0, 1\}, \quad \forall t, \quad (1.4j)$$

$$y_{jtt^ut^a}^O, r_{jtt^u}, Q_{jtt^u}, K_{tt^a}^O \in \mathbb{Z}^+, \quad \forall j, t, t^a, t^u \leq t. \quad (1.4k)$$

To model whether or not demand exceeds capacity, we introduce binary variables  $z_t^I$  and  $z_t^E$ . We let  $z_t^I = 1$  if  $D_t^I \leq \alpha^I K$  at time  $t$  and  $z_t^I = 0$  otherwise. Similarly,  $z_t^E = 1$  if  $D_t^E \leq \alpha^E K$  and

0 otherwise. Constraints (1.4c)-(1.4f) ensure that  $z_t^I$  and  $z_t^E$  take on the correct values depending on whether or not demand is less than the corresponding capacity. Constraint set (1.4g) ensures that the total number of next day appointments given to outpatients does not exceed the remaining capacity for outpatients plus the unused capacity that was reserved for inpatients. Constraint set (1.4h) is similar to (1.4g), *i.e.*, it ensures that the total number of outpatient appointments does not exceed the remaining capacity on the subsequent days. The only difference is that in (1.4g) we also have the unused capacity that was initially allocated for inpatients which can now be used for outpatients. Constraint set (1.4i) ensures that patients do not wait more than  $T^u$  days in the queue. Finally, constraints (1.4j) and (1.4k) are the binary and integrality constraints.

Note that, when solving the first-stage problem (2SIP), the objective function (1.4a) and the constraint set (1.4g) are nonlinear. However, we will not linearize these since our approximation and decomposition approaches will not require solving (2SIP) directly. Instead, we will reformulate the problem as described in Section 1.4.

**Limitations of the model:** One of the limitations of our TSSP model is that is anticipative, *i.e.*, it relies on knowing the demand for the whole planning horizon. Another limitation is that the model assumes the system is initially empty. Also, the model is considering a finite planning horizon which may lead to end-of-horizon effects. We address these limitations to some extent as discussed later in Sections 1.4 and 1.5.

## 1.4 Solution Approach

Due to the curse of dimensionality, solving (2SIP) as presented in Section 1.3 is impractical. To overcome this complexity, we develop a sample average approximation (SAA) approach to generate tight upper and lower bounds. The SAA procedure generates a random sample  $\omega^1, \omega^2, \dots, \omega^S$  of  $S$  scenarios from  $\Omega$ , the set of all possible scenarios, and solves  $M$  replications of the following deterministic SAA problem:

$$(2\hat{\text{SIP}}) \quad \hat{C}_S = \min_{\alpha} \quad \frac{1}{S} \sum_{s=1}^S C(\alpha, \omega^s) \tag{1.5a}$$

$$\text{s.t.} \quad (1.3b), (1.3c). \tag{1.5b}$$

Note that  $\hat{C}_S \rightarrow C^*$  as  $S \rightarrow \infty$ , and estimates of the optimal first-stage solutions for the original stochastic problem can be obtained by solving this deterministic version [88]. Algorithm 1 below formalizes our proposed SAA approach. As shown in the algorithm, the average of the  $M$  replications ( $\bar{C}_S$ ) provides a statistical lower bound for  $C^*$ . For each solution to (2 $\hat{\text{SIP}}$ ) from the  $M$  replications, the second-stage problem (1.4) is solved using a larger sample size  $S'$ . Among this larger sample, the one with the smallest objective value ( $\hat{C}_{S'}(\hat{\alpha}^*)$ ) is our statistical upper bound for  $C^*$ . We also calculate the variances of the lower and upper bound estimates, *i.e.*,  $\sigma_{\bar{C}_S}^2$  and  $\sigma_{\hat{C}_{S'}(\hat{\alpha}^*)}^2$ , respectively. The proofs of the estimation of lower and upper bounds are provided by [60] and [88], and thus, omitted here. The algorithm increases the sample sizes  $S$  and  $S'$  until the optimality gap and the variance of the gap estimator are small.

---

**Algorithm 1** Sample Average Approximation (SAA)

---

**Step 1:** Initialize  $S$ ,  $S'$ , and  $M$ ;  
**Step 2:** **For**  $m = 1, \dots, M$   
    Solve (2 $\hat{\text{SIP}}$ ) using DBB to obtain objective values  $\hat{C}_S^m$  and solutions  $\hat{\alpha}^m$ ;  
**Step 3:** Calculate  $\bar{C}_S = \frac{1}{M} \sum_{m=1}^M \hat{C}_S^m$  and  $\sigma_{\bar{C}_S}^2 = \frac{1}{M(M-1)} \sum_{m=1}^M (\hat{C}_S^m - \bar{C}_S)^2$ ;  
**Step 4:** **For** each  $\hat{\alpha}^m$   
    Solve (1.4) and compute  $\hat{C}_{S'} = \frac{1}{S'} \sum_{s=1}^{S'} C(\hat{\alpha}^m, \omega^s)$  and  $\sigma_{\hat{C}_{S'}(\hat{\alpha})}^2 = \frac{1}{S'(S'-1)} \sum_{s=1}^{S'} (C(\hat{\alpha}^m, \omega^s) - \hat{C}_{S'})^2$ ;  
**Step 5:** Let  $\hat{\alpha}^* = \arg \min \left\{ \hat{C}_{S'}(\hat{\alpha}) : \hat{\alpha} \in \{\hat{\alpha}^1, \dots, \hat{\alpha}^M\} \right\}$ ;  
**Step 6:** Calculate  $\Delta_C = \hat{C}_{S'}(\hat{\alpha}^*) - \bar{C}_S$  and  $\sigma^2 = \sigma_{\bar{C}_S}^2 + \sigma_{\hat{C}_{S'}(\hat{\alpha}^*)}^2$ ;  
**Step 7:** **If** ( $\Delta_C < \epsilon$  and  $\sigma^2 < \epsilon$ ) **then** report  $\hat{\alpha}^*$  as the optimal solution and terminate;  
    **Else** increase  $S$  and  $S'$  and go back to Step 2.

---

Solving (2 $\hat{\text{SIP}}$ ) in Algorithm 1, while easier than solving (2SIP), is still a challenging task for large  $S$ . To that end, we developed a decomposition based branch-and-bound (DBB) algorithm which was originally proposed by [7] to solve TSSP models with continuous first-stage and discrete second-stage variables. The main idea behind DBB is to partition the search space to efficiently identify candidate solutions [6, 7]. To be able to implement DBB and ensure convergence, the following assumptions must be satisfied (all of which are satisfied for (2 $\hat{\text{SIP}}$ )): (A1) The uncertain parameter  $\omega$  follows a discrete distribution with finite support. (A2) The first-stage constraint set is nonempty and compact. (A3) The second-stage variables are purely integer. (A4) The technology matrix is deterministic. (A5) For each scenario the second-stage problem is bounded. (A6) For each

scenario, the second-stage constraint matrix is integral. We reformulate (2 $\hat{\text{SIP}}$ ) as follows:

$$(\text{TP}) \quad \min_{\boldsymbol{\chi}} \quad f(\boldsymbol{\chi}) \tag{1.6a}$$

$$\text{s.t.} \quad \boldsymbol{\chi} \in X, \tag{1.6b}$$

where  $f(\boldsymbol{\chi}) = \frac{1}{S} \sum_{s=1}^S \Psi^s(\boldsymbol{\chi})$ ,  $\Psi^s(\boldsymbol{\chi}) = \min\{f^s \mathbf{y} \mid D^s \mathbf{y} \geq h^s + \boldsymbol{\chi}, \mathbf{y} \in Y \cap \mathbb{Z}\}$ , and  $X = \{\boldsymbol{\chi} \mid \boldsymbol{\chi} = T\boldsymbol{\alpha}, \text{ with (1.3b) and (1.3c)}\}$ . In the stochastic programming literature, the matrix  $T$  is known as the *technology matrix* and variables  $\boldsymbol{\chi}$  as the *tender variables* that link the first- and second-stage problems. Note that for our problem  $T$  is deterministic, *i.e.*, it is independent of the scenario observed. The term  $\Psi^s(\boldsymbol{\chi})$  is essentially a compact representation of the second-stage problem given by formulation (1.4) where  $\mathbf{y}$  represents the collection of all second-stage decision variables (*i.e.*,  $\mathbf{y} = (\mathbf{y}^O, \mathbf{r}, \mathbf{Q}, \mathbf{K}^O)$ ),  $f^s$  represents the objective function (1.4a), and  $D^s$ ,  $h^s$ , and  $Y$  represent the constraints (1.4b)-(1.4k) with  $D^s$  corresponding to the scenario dependent coefficients,  $h^s$  the scenario dependent constants,  $Y$  the scenario independent constraints, and  $T$  the scenario independent parts of the constraint set which include the first-stage variables. This reformulation allows us to consider a larger number of scenarios in Step 2 of Algorithm 1. More specifically, the DBB algorithm below enables us to avoid solving (2 $\hat{\text{SIP}}$ ) directly. Instead of the first-stage variables, we search the space of the tender variables for global optima. The search space of  $\boldsymbol{\chi}$  is partitioned into subsets of the form  $\prod_j (l_j, u_j]$ , for each component  $j$  of  $\boldsymbol{\chi}$  where  $l_j$  is a point at which the second-stage value function ( $\Psi^s(\cdot)$ ) may be discontinuous [7]. By branching this way, we isolate subsets over which the second-stage value function is constant. Thus, we can solve (2 $\hat{\text{SIP}}$ ) exactly.

---

**Algorithm 2** Decomposition based Branch-and-Bound (DBB)

---

- Step 1:** Initialize  $U = \infty$ ,  $k = 0$ ,  $\mathcal{P}^k$ , and  $\mathcal{L}$ ;  
**Step 2:** **If** ( $\mathcal{L} = \emptyset$ ) **then** terminate with solution  $\hat{\boldsymbol{\chi}}^*$ ;  
    **Else** select and remove a subproblem  $k$  from  $\mathcal{L}$  (*i.e.*,  $\mathcal{L} = \mathcal{L} \setminus \{k\}$ );  
**Step 3:** Generate upper ( $\gamma^k$ ) and lower ( $\beta^k$ ) bounds for subproblem  $k$ ;  
**Step 4:** Set  $L = \min_{i \in \mathcal{L} \cup \{k\}} \beta^i$ ;  
**Step 5:** **If** ( $\gamma^k < U$ ) **then** set  $U = \gamma^k$  and  $\boldsymbol{\chi}^* = \boldsymbol{\chi}^k$ ;  
**Step 6:** Fathom the subproblem (*i.e.*, set  $\mathcal{L} = \mathcal{L} \setminus \{i \mid \beta^i > U\}$ );  
**Step 7:** **If** ( $\beta^k > U$ ) **then** go to Step 2;  
**Step 8:** Branch by partitioning  $\mathcal{P}^k$  into  $\mathcal{P}^{k_1}$  and  $\mathcal{P}^{k_2}$ ;  
**Step 9:** Set  $\mathcal{L} = \mathcal{L} \cup \{k_1, k_2\}$ ,  $\beta^{k_1} = \beta^k$ ,  $\beta^{k_2} = \beta^k$ ,  $k = k + 1$ , and go to Step 2.
- 

In Step 1 of Algorithm 2, we begin (after setting  $k = 0$ ) by constructing the hyper-rectangle

$\mathcal{P}^k = \prod_j (l_j^k, u_j^k] \supset X$  and adding the problem  $\inf\{f(\chi) | \chi \in X \cap \mathcal{P}^k\}$  to the list of open subproblems  $\mathcal{L}$ . For each component  $j$  of  $\chi$ , we set  $l_j = \min\{\chi_j | \chi \in X\}$  and  $u_j = \max\{\chi_j | \chi \in X\}$  where the optimization problems are linear programs since  $X$  is polyhedral. Then, for each  $j$  and scenario  $s$ , we find  $k_j^s \in \mathbb{Z}$  such that  $k_j^s - h_j^s - 1 < l_j < k_j^s - h_j^s$ . If  $l_j + h_j^s$  is integral then set  $k_j^s = l_j + h_j^s$ ; otherwise  $k_j^s = \lfloor l_j + h_j^s + 1 \rfloor$ . Finally, we set  $l_j^k = \max_s \{k_j^s - h_j^s - 1\}$  and  $u_j^k = u_j$ . In Step 2, to ensure convergence of the algorithm (proof given by [7]), we select subproblem  $k$  such that  $\beta^k = L$ . In Step 3, for a given subset  $\mathcal{P}^k$ , we obtain a lower bound on the corresponding subproblem by solving the following formulation:

$$\text{(LB)} \quad \beta^k = \min \quad \theta \tag{1.7a}$$

$$\text{s.t.} \quad (1.6b), \tag{1.7b}$$

$$l^k \leq \chi \leq u^k, \tag{1.7c}$$

$$\theta \geq \sum_{s=1}^S \frac{1}{S} \Psi^s(l^k + \epsilon). \tag{1.7d}$$

In problem (LB),  $\Psi^s(\cdot)$  is constant over  $(l^k, l^k + \epsilon]$  for all  $s$  when  $\epsilon$  is sufficiently small [7]. The value of  $\epsilon$  can be calculated *a priori* using the following algorithm:

---

**Algorithm 3** Calculation of  $\epsilon$

---

**Step 1:** For each component  $j$  of  $\chi$

Set  $s = 1, \Xi = \emptyset$ . Choose  $k_j^1 \in \mathbb{Z}$ . Let  $\chi_j^0 = k_j^1 - h_j^1 - 1$  and  $\chi_j^1 = \chi_j^0 + 1$ . Set  $\Xi = \Xi \cup \{\chi_j^0, \chi_j^1\}$ ;

**Step 2:** For  $s = 2, \dots, S$

Set  $k_j^s = \lfloor \chi_j^1 + h_j^s \rfloor$ . Let  $\chi_j^s = k_j^s - h_j^s$ ;

If  $\Xi \cap \{\chi_j^s\} = \emptyset$  then set  $\Xi = \Xi \cup \{\chi_j^s\}$ ;

**Step 3:** Sort the elements of  $\Xi$  such that  $\chi_j^0 = \xi_j^0 < \xi_j^1 < \dots < \xi_j^n = \chi_j^1$  with  $n \leq S$ ;

**Step 4:** Let  $\epsilon_j = \min_i \{\xi_j^i - \xi_j^{i-1}\}$ ;

**Step 5:** Set  $\epsilon = \frac{1}{2} \min_j \{\epsilon_j\}$ .

---

In Step 3 of Algorithm 2, we also generate an upper bound. For a given subset  $\mathcal{P}^k$  such that  $\mathcal{P}^k \cap X \neq \emptyset$ , let  $\chi^k$  be an optimal solution to (LB). Since  $\chi^k$  is feasible to (TP) we can simply set  $\gamma^k = f(\chi^k)$ . Finally, in Step 8 we perform branching. To do this we identify the variable  $j'$  by determining the value of  $\chi_{j'}$  where the the current second-stage problem becomes infeasible. For each scenario  $s$ , let  $y^s$  be the solution of the second-stage subproblem when solving (LB). Then, for each  $j$  compute  $p_j = \min_s \{(D^s y^s)_j - h_j^s\}$ . Let  $j' \in \operatorname{argmax}_j \{\min\{p_j - l_j^k, u_j^k - p_j\}\}$  and split  $\mathcal{P}^k$



into  $\mathcal{P}^{k_1} = (l_{j'}^k, p_{j'}] \prod_{j \neq j'} (l_j^k, u_j^k]$  and  $\mathcal{P}^{k_2} = (p_{j'}, u_{j'}^k] \prod_{j \neq j'} (l_j^k, u_j^k]$ .

## 1.5 Numerical Study

To evaluate the advantages of postponement in making acceptance and scheduling decisions about outpatient appointments, we conducted an extensive numerical study. We also performed a sensitivity analysis to demonstrate how the performance is affected by changes in some problem parameters.

### 1.5.1 Input data

The patient arrival rates  $\lambda_j, \lambda^I, \lambda^E$  and the parameters in Table 1.1 are required input for the proposed model. We consider two priority classes of outpatients ( $j = 1, 2$ ). The values of  $\lambda_j, \lambda^I, \lambda^E$  are estimated based on the average arrival rates of different patient types at the Radiology Department of Prisma Health. The parameters in Table 1.1 are estimated with the assistance of physicians at the Radiology Department. However, due to confidentiality concerns we only present normalized values in the appendix. The daily regular capacity of the clinic is estimated to be  $K = 175$  appointments, and the planning horizon is set to  $T = 50$  days. In our analysis we ignored the first week (*i.e.*, we used it as our warm-up period). We also ignored the last week of the planning horizon to eliminate any end-of-horizon effects.

### 1.5.2 Experimental setup

The problem is implemented in C<sup>++</sup>. The decomposed problems are solved on an Intel Core i7-9700 CPU utilizing the Gurobi 7.0 solver. The computational time required to implement the SAA algorithm will grow as  $S$ ,  $S'$ , and  $M$  increase in Algorithm 1. The growth can be linear or exponential depending on whether or not a decomposition approach is used [51]. In our case the growth is linear since we are using DBB as presented in Algorithm 2. Our first set of experiments were conducted to determine suitable values for  $S$ ,  $S'$ , and  $M$ . We began with initializing  $S=10$  and  $S'=100$ , and increased these values in increments of 10 until  $\Delta_c$  and  $\sigma^2$  values were less than  $\epsilon = 0.01$ . We also tested different values for  $M$  from the set  $\{10, 20, 30, 40, 50\}$ . The final values for  $S$ ,  $S'$ , and  $M$  were respectively, 100, 500, and 30.

### 1.5.3 The base scenario

After determining the values for  $S$ ,  $S'$ , and  $M$ , we conducted a large number of experiments to test and compare the performance of four different appointment scheduling policies.

#### 1.5.3.1 Policy 1:

This policy refers to what is currently being used by the clinic. As described earlier, the clinic currently allocates a portion of the regular capacity for emergency patients, a portion to inpatients, and uses the remaining capacity for outpatients. Capacity allocations are done based on the  $\alpha^I$  and  $\alpha^E$  values obtained from the optimization model without postponement (which is explained below in policy 3). In policy 1, appointment decisions are made as soon as a patient arrives on an FCFS basis and the capacities are dedicated. When an emergency patient arrives that patient is seen immediately if there is capacity, otherwise they are referred to another clinic. When an inpatient arrives that patient is given the earliest next day appointment as long as there is capacity, otherwise they are seen during overtime hours. When an outpatient arrives they are given the earliest possible appointment (regardless of type) over the next seven days. If there is no capacity then they are referred to another clinic.

#### 1.5.3.2 Policy 2:

This is the proposed policy where the outpatients are kept in an acceptance queue up to 72 hours (3 days). The emergency patients are handled the same way as in policy 1, but the appointment decisions for inpatients and outpatients are made at the end of each day. The SAA and DBB approaches are used in policy 2 to determine the capacity allocations in stage 1 and appointment decisions in stage 2. Note that policy 2 is anticipative, *i.e.*, the demand for the whole planning horizon is revealed at the beginning of stage 2. One can think of policy 2 as the policy with perfect information and postponement.

#### 1.5.3.3 Policy 3:

This policy is similar to policy 1 in that outpatients are not kept in an acceptance queue. On the other hand, policy 3 is similar to policy 2 because it is anticipative and uses the same SAA and DBB approaches to make capacity allocation and appointment scheduling decisions. The main

difference is that the regular working hours of a day is split into  $T' = 54$  periods. In other words, in policy 3 appointment decisions are made every 10 minutes, *i.e.*, near real-time. The decisions regarding acceptance and referral of outpatients are taken in each decision epoch  $t'$  ( $t' = 1, 2, \dots, T'$ ). The original arrival rates are divided by  $T'$  and constraint set (1.4g) is modified to reflect the fact that unused capacity that is allocated for inpatients cannot be used in policy 3. Additionally, constraint sets (1.1a) and (1.1b), which capture the evolution of the acceptance queue, are removed from the model. Table 2 in the Appendix shows the values of the parameters for policy 3. Policy 3 is essentially the same as policy 1 but with perfect information.

#### 1.5.3.4 Policy 4:

Based on our observations of the optimal solutions from policy 2, we developed a simple benchmark policy, which is non-anticipative (*i.e.*, does not rely on knowing the demand for the whole planning horizon) and does not keep patients in an acceptance queue. In policy 4 acceptance or referral of all patients are done on arrival but in a way that mimics the decisions made under policy 2. The optimal values for  $\alpha^I$  and  $\alpha^E$  obtained from policy 2 are used for capacity allocation. In policy 4 the emergency patients and inpatients are handled the same way as in policy 1. The outpatients, on the other hand, are handled differently. In policy 1 all outpatients are given earliest available appointments on arrival. In policy 4, however, some of the outpatients are referred to other clinics on arrival regardless of available capacity. As will be shown later, 87.37% of outpatients that request an appointment end up getting one in policy 2. More specifically, the average acceptance rates are 85% and 92.5%, respectively, for Type 1 and Type 2 outpatients. Thus, 15% (7.5%) of Type 1 (Type 2) outpatients are immediately referred to another clinic on arrival in policy 4. For those outpatients who are not referred to another clinic, Type 2 outpatients are given the earliest available appointment beginning with day two of the planning period. In other words, next day appointments are not given to Type 2 outpatients. For Type 1 outpatients decisions are made based on  $b_{jtu^a}$  values. Since policy 2 keeps these patients in the acceptance queue for two days, the  $b_{1,2,t^a}$  values are sorted in non-decreasing order, and appointments are given based on this order. However, if the cost difference in consecutive days are within 20% of each other then the later day in the horizon is selected. For the clinic in our case study, this translates to considering days two, five, and seven from time of arrival for possible appointments. We begin with day two and check the remaining capacity. If this remaining capacity is more than 33% of the total daily outpatient capacity  $K(1 - \alpha_I - \alpha_E)$

then the Type 1 outpatient is given an appointment on that day with probability 0.75 (based on our observation of policy 2). For a Type 1 outpatient that does not get an appointment on day two the next option (*i.e.*, day five) is considered and the same rules are applied. This process is repeated until the last day of the planning horizon (day seven in this particular example) at which point all remaining Type 1 outpatients are given an appointment on this last day.

Table 1.3 compares policy 1 to policy 2. As can be see from the table, policy 2 significantly improves the expected average cost for the clinic. The improvement ranged from about 40% to 46% depending on how long outpatients are allowed to be kept in the acceptance queue. Recall that the clinic is not willing to keep the outpatients in the acceptance queue more than 3 days. Based on our experiments, the lowest total cost was achieved when  $T^u = 2$ . Thus, in our remaining experiments the value of  $T^u$  is fixed at 2.

Table 1.3: Cost improvement and capacity allocation for the base scenario

|                       | Policy 1 | Policy 2  |           |           |
|-----------------------|----------|-----------|-----------|-----------|
|                       |          | $T^u = 1$ | $T^u = 2$ | $T^u = 3$ |
| Avg. cost improvement | -        | 40.7%     | 45.5%     | 39.9%     |
| $\alpha_I$            | 13%      | 20%       | 20%       | 20%       |
| $\alpha_E$            | 34%      | 34%       | 34%       | 35%       |

Table 1.4 summarizes the results for all four policies. The main takeaway is that policy 4, which is very easy to implement, performs really well. Recall that policy 1 is the current policy used at the clinic, policy 3 is the “optimal” version of policy 1, policy 2 is the one that keeps outpatients in an acceptance queue, and policy 4 is a simple heuristic that we developed which mimics policy 2. As can be seen from Table 1.4, policies 2 and 4 result in more than 45% cost improvement compared to policy 1. The small cost difference between policies 1 and 3 (and between policies 2 and 4) suggest that the value of perfect information is minimal. While this may sound counter-intuitive it is expected, because policies 2 and 3 are run using large samples of scenarios and the average is reported. With respect to emergency patients almost none are referred to another clinic under all four policies. With respect to inpatients policies 2 and 4 handle almost all of them during regular hours, but policies 1 and 3 handle about 3.5% of them during overtime hours.

The main difference among the four policies is in the way they handle outpatients. In

policy 1, about 82% of out patients are handled during regular hours and the rest are referred to other clinics. Note that this percentage stays almost the same for Type 1 and Type 2 outpatients which makes sense because the current policy functions on an FCFS basis and does not prioritize outpatients. On the other hand, in policy 3 almost all of Type 2 outpatients are seen during regular hours but only 75% of Type 1 outpatients are seen during regular hours with an overall average of almost 84% for all outpatients. Recall that policy 3 is the anticipative version of policy 1. Thus, knowing the demand for the whole planning horizon allows policy 3 to prioritize different outpatients. In policy 2 the percentage of outpatients seen during regular hours is more than 87%, a relatively significant increase over the current system. Policy 2 is able to do this because it is able to utilize the unused capacity allocated to inpatients. For a fair comparison, the outpatients acceptance percentages for policy 2 include those patients that leave the acceptance queue. For example, if 100 Type 1 outpatients arrive then about 5 leave the acceptance queue and of the remaining 95, on average, 85 get appointments and 10 are referred to other clinics. The performance of policy 4 with respect to patient acceptance is very similar to policy 2 since it was designed to mimic policy 2.

Table 1.4: Summary of all results for all four policies for the base scenario

|                                      | Policy 1 | Policy 2 | Policy 3 | Policy 4 |
|--------------------------------------|----------|----------|----------|----------|
| Avg. cost improvement                | -        | 45.5%    | 2.5%     | 45.2%    |
| Emergency patient acceptance         | 99.6%    | 99.6%    | 99.6%    | 99.6%    |
| Inpatient acceptance                 | 96.4%    | 99.9%    | 96.5%    | 99.9%    |
| Outpatient acceptance                | 81.7%    | 87.4%    | 83.9%    | 86.5%    |
| Type 1 outpatient acceptance         | 81.5%    | 85.0%    | 75.0%    | 84.2%    |
| Type 2 outpatient acceptance         | 82.0%    | 92.5%    | 99.9%    | 92.5%    |
| Type 1 outpatients leaving the queue | -        | 5.0%     | -        | -        |
| Type 2 outpatients leaving the queue | -        | 5.0%     | -        | -        |
| Inpatient capacity not used          | 4.4%     | 0.0%     | 4.5%     | 0.0%     |
| Days in acceptance queue (Type 1)    | -        | 1.9      | -        | -        |
| Days in acceptance queue (Type 2)    | -        | 1.0      | -        | -        |
| Appointment days ahead (Type 1)      | 6.2      | 5.5      | 6.4      | 5.6      |
| Appointment days ahead (Type 2)      | 6.1      | 3.1      | 6.0      | 3.1      |
| $\alpha^I$                           | 13%      | 20%      | 13%      | 20%      |
| $\alpha^E$                           | 34%      | 34%      | 34%      | 34%      |
| Solution time (sec.)                 | 0.10     | 120      | 110      | 0.14     |

Table 1.4 also shows the percentage of the overall capacity allocated to each patient group.

Recall that policy 1 (policy 4) simply uses the  $\alpha^I$  and  $\alpha^E$  values obtained from policy 3 (policy 2). Both of policies 2 and 3 allocate about 34% of the capacity to emergency patients. However, policy 2 allocates more capacity to inpatients compared to policy 3. While the increase from 13% to 20% may seem unnecessary, it is expected because under policy 2 with postponement the extra capacity allocated for inpatients can be used for outpatients when needed.

Another interesting observation is related to the indirect waiting times. As seen in Table 1.4, under the current policy, the indirect waiting times for Type 1 and Type 2 outpatients are 6.2 and 6.1 days, respectively. As expected, policy 1 does not distinguish between the two types of outpatients. In policy 3, since it is anticipative, the indirect waiting times are 6.4 and 6.0 favoring Type 2 outpatients slightly, but the overall average is almost the same as in policy 1. In policy 2 the indirect waiting times are 7.4 (1.9+5.5) and 4.1 days, respectively, for Type 1 and Type 2 outpatients. This shows that policy 2 prioritizes Type 2 outpatients. The average waiting time is decreased by about 2 days for Type 2 outpatients in the expense of about a 1 day increase for Type 1 outpatients. Policy 4 mimics policy 2 but it does not keep an acceptance queue. Thus, policy 4 does very well in reducing indirect waiting times.

With respect to computational time, policies 1 and 4 are very fast since they are essentially simulating the appointment system using simple rules. In addition, policies 1 and 4 do not compute capacity allocations but use the values obtained from policies 3 and 2, respectively. Thus, the average CPU time per replication is 0.10 and 0.14 seconds, respectively, for policies 1 and 4. Policies 2 and 3 solve complicated optimization problems to optimize capacity allocations and appointment schedules. As such the average CPU time per replication is 120 and 110 seconds for policies 2 and 3, respectively.

#### 1.5.4 Sensitivity analysis

The results presented in Section 1.5.3 demonstrate the effectiveness of policy 2 on the base scenario (referred to as experiment 1). While policy 4 is our proposed policy (because it is easy to implement), it is based on policy 2. Thus, in this section, we perform a sensitivity analysis to observe how policy 2 performs under different conditions. For this analysis, the values of the following parameters are changed one at a time:  $b_{jt^u t^a}$ ,  $c^I$ ,  $c^E$ ,  $\lambda^I$ , and  $p_{jt^u}$ .

#### 1.5.4.1 Scheduling costs of outpatients:

To understand the effect of changing outpatient scheduling costs on the optimal solution, we increased  $b_{jt^u t^a}$  by 50% and 100% in experiments 2 and 3, respectively. By increasing all the  $b_{jt^u t^a}$  values with the same percentage we penalize both wait times in the acceptance queue and the time until appointments in experiments 2 and 3. In experiments 4 and 5, we penalize long wait times in the acceptance queue by increasing the  $b_{jt^u t^a}$  values for only  $t^u = 2$  by 50% and 100%, respectively. In experiments 6 and 7, we penalize scheduling later appointments, where  $b_{jt^u t^a}$  values for  $t^a \geq 3$  are increased by 50% and 100%, respectively. Figures 1.1 and 1.2 represent the results of experiments 1-7. As seen in Figure 1.1b, the percent of capacity allocated for emergency patients ( $\alpha^E = 0.34$ ) is not impacted by changes to  $b_{jt^u t^a}$ . On the other hand, capacity allocated for inpatients ( $\alpha^I$ ) increases slightly from 20% to 22% in experiment 3 since we are willing to reserve more next day appointments for outpatients. The box plots in Figure 1.1a show the normalized values of average total cost for the system. In other words, the average total cost for experiment 1 is normalized to 100, and thus, the other values show the corresponding change in cost. As expected, the total cost increases the most in experiments 2 and 3 since all  $b_{jt^u t^a}$  values are increased here whereas only a subset of the  $b_{jt^u t^a}$  values are increased in experiments 4-7.

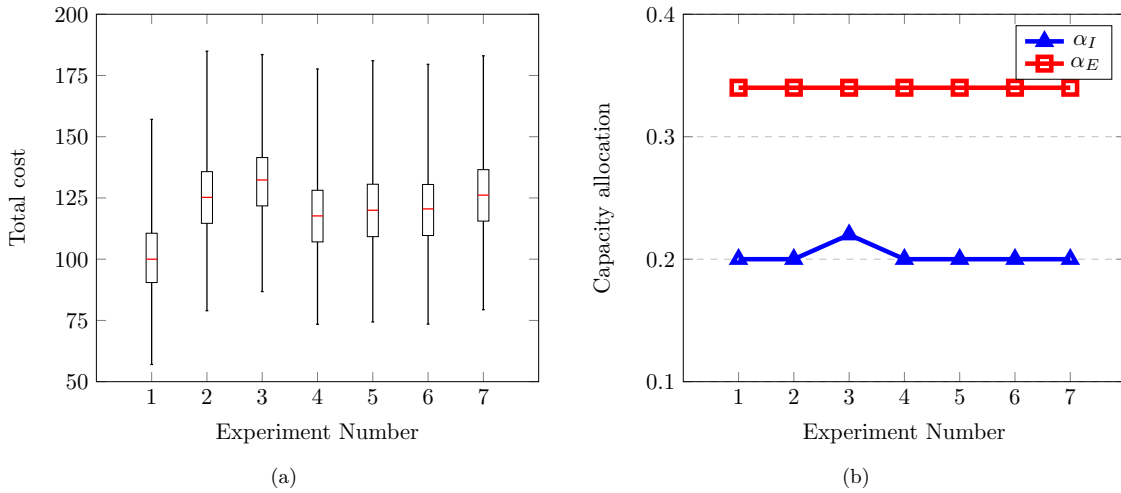


Figure 1.1: Effect of changing  $b_{jt^u t^a}$  on average total cost and capacity allocation

Figure 1.2a shows that in the base scenario, all Type 2 outpatients wait in the acceptance queue for one day before they are given an appointment. Type 1 outpatients, however, wait for almost two days in the queue. As the cost of waiting in the queue increases the waiting time

for Type 2 outpatients remain the same. For Type 1 outpatients, it decreases. In other words, acceptance and referral decisions are made sooner. Figure 1.2b shows that time from the day of acceptance to the day of appointment decreases as  $b_{jt^u t^a}$  increases. Because all  $b_{jt^u t^a}$  are increased in experiments 2 and 3, the indirect waiting time decreases sharply. In experiments 6 and 7, the  $b_{jt^u t^a}$  values were increased for only high values of  $t^a$ , as such, compared to the base case the drop in indirect waiting time is not as dramatic. However, policy 2 tries to offer earlier appointments to both outpatient types in experiments 6 and 7. Additional insight on these experiments are also discussed in Section 1.5.4.5 based on Tables 1.5 and 1.6.

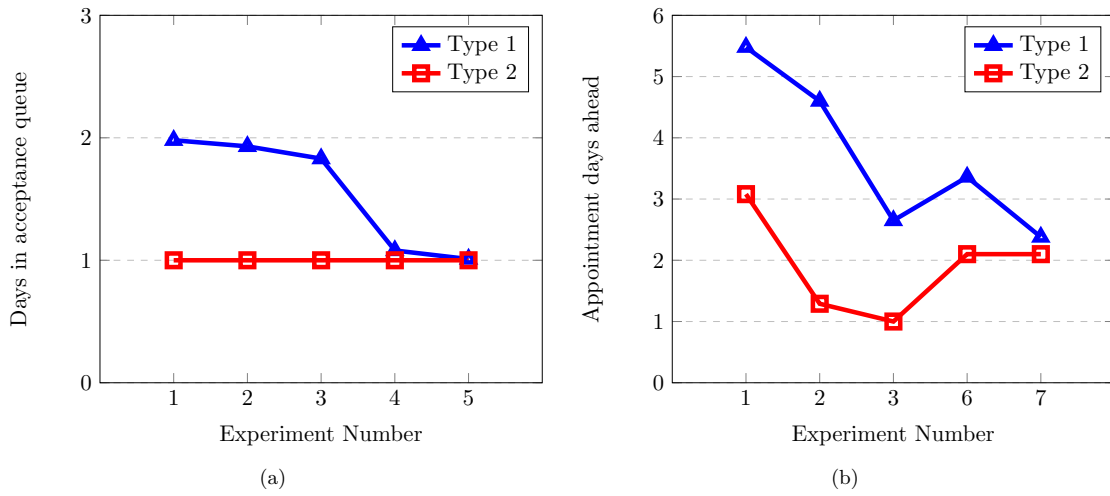


Figure 1.2: Effect of changing  $b_{jt^u t^a}$  on the number of days in acceptance queue and indirect waiting time

#### 1.5.4.2 Referral and overtime costs of emergency patients and inpatients:

Emergency patients are the highest priority patients followed by inpatients. Lack of available capacity to schedule them during regular hours results in additional cost, specifically, overtime cost  $c^I$  for inpatients and referral cost  $c^E$  for emergency patients. Experiments 8 and 9 measure the effect of increasing these parameters by 50% and 100%, respectively. With respect to capacity allocation,  $\alpha^E$  and  $\alpha^I$  remained the same in experiments 8 and 9 as they were in experiment 1. The only difference was on the average total cost, which increased by about 4% from experiment 1 to 8 and about 8% from experiment 1 to 9.



#### 1.5.4.3 Inpatient arrivals:

Since capacity allocated to inpatients can also be used for outpatients in policy 2, arrival rate of inpatients affect the scheduling of outpatients. To observe this impact, we performed experiments 10 and 11 where  $\lambda^I$  is increased by 25% and 50%, respectively. As one would expect, increasing the inpatient arrival rate resulted in higher  $\alpha^I$  values (25% in experiment 10 and 28% in experiment 11). On the other hand, there was no change to the capacity allocated to emergency patients which stayed at 34%. However, the outpatient acceptance rate decreased in both of these experiments. Since more inpatients are arriving into the system, there is less capacity left for outpatients and more of them are referred to other clinics. Given that all problem parameters (including total capacity) remained the same as  $\alpha^I$  was increased, the total system cost increased about 28% and 44% in experiments 10 and 11, respectively, compared to experiment 1.

#### 1.5.4.4 Abandonment rate of outpatients:

As discussed earlier, one of the disadvantages of implementing postponed acceptance in a service system such as a diagnostic clinic is that customers may abandon the acceptance queue. We assume that a proportion  $(1 - p_{jt^u})$  of type  $j$  outpatients leave the queue after having waited  $t^u$  days. In addition to the outpatient type and the amount of time they have waited in the queue, the abandonment rates also depend on the type of diagnostic clinic. To capture the effect of such changes, we decreased the value of  $p_{jt^u}$  by 5% in experiment 12 and 10% in experiment 13. In other words, the chances of an outpatient abandoning the acceptance queue is higher in experiments 12 and 13. As a result,  $\alpha^I$  stayed at 20% in experiment 12 but slightly increased to 21% in experiment 13. Since fewer outpatients are waiting for an appointment due to abandonment, the rejection rate of outpatients decreased in experiments 12 and 13 compared to experiment 1. More specifically, in experiment 12 the number of outpatients referred to other clinics decreased by 85%. In experiment 13 no outpatient was referred to another clinic. This resulted in lower average total cost with a decrease of about 27% and 40% in experiments 12 and 13, respectively, as compared to experiment 1.

#### 1.5.4.5 Additional insights:

Table 1.5 provides the percentage of outpatients who have received appointments after waiting one day or two days in the acceptance queue. Type 2 outpatients always received an appointment after only one day in the acceptance queue. On the other hand, majority of type 1 outpatients wait for two days in the acceptance queue in all of the experiments except experiments 4 and 5. Recall that in experiments 4 and 5 the  $b_{jt^u t^a}$  values are increased for  $t^u = 2$  by 50% and 100%, respectively. In other words, waiting in the acceptance queue for two days is costly in these cases. Thus, in experiments 4 and 5 most of the outpatients get an appointment after only one day in the queue.

Table 1.5: Percentage of outpatients waiting in the acceptance queue for one vs. two days

| Experiment | Type 1 outpatients<br>One day in queue | Type 1 outpatients<br>Two days in queue | Type 2 outpatients<br>One day in queue | Type 2 outpatients<br>Two days in queue |
|------------|--|---|--|---|
| 1          | 1.93                                   | 98.07                                   | 100.00                                 | 0.00                                    |
| 2          | 6.63                                   | 93.37                                   | 100.00                                 | 0.00                                    |
| 3          | 16.60                                  | 83.40                                   | 100.00                                 | 0.00                                    |
| 4          | 92.14                                  | 7.86                                    | 100.00                                 | 0.00                                    |
| 5          | 99.40                                  | 0.60                                    | 100.00                                 | 0.00                                    |
| 6          | 7.61                                   | 92.39                                   | 100.00                                 | 0.00                                    |
| 7          | 1.97                                   | 98.03                                   | 100.00                                 | 0.00                                    |
| 8          | 1.70                                   | 98.44                                   | 100.00                                 | 0.00                                    |
| 9          | 1.56                                   | 98.43                                   | 100.00                                 | 0.00                                    |
| 10         | 0.54                                   | 99.46                                   | 100.00                                 | 0.00                                    |
| 11         | 0.22                                   | 98.78                                   | 100.00                                 | 0.00                                    |
| 12         | 2.18                                   | 97.82                                   | 100.00                                 | 0.00                                    |
| 13         | 3.59                                   | 96.41                                   | 100.00                                 | 0.00                                    |

The unused inpatient capacity before and after scheduling outpatients out of the acceptance queue are captured and listed in columns 2 and 3 of Table 1.6, respectively. Columns 4, 5, and 6 of the table shows the percentage of the patients that ultimately received appointments. Depending on the problem parameters, the unused inpatient capacity varies between 44% and 56%. Note that almost all of the inpatients receive appointments and the leftover capacity is used for outpatients. Outpatient acceptance rate is low for experiment 11 but over 80% in all of the other experiments. Recall that in experiments 10 and 11 the inpatient arrival rate was increased, as such, there is not much leftover capacity that can be used for outpatients compared to the other experiments.

Table 1.6: Capacity utilization of each patient type for policy 2

| Experiment | Unused inpatient capacity<br>before scheduling<br>outpatients (%) | Unused inpatient capacity<br>after scheduling<br>outpatients (%) | Inpatient<br>acceptance (%) | Emergency patient<br>acceptance (%) | Outpatient<br>acceptance (%) |
|------------|---|--|-----------------------------|-------------------------------------|------------------------------|
| 1          | 44.45   | 0.00   | 99.99                       | 99.56                               | 87.39                        |
| 2          | 44.45   | 0.00   | 99.99                       | 99.56                               | 85.15                        |
| 3          | 52.78   | 0.00   | 99.99                       | 99.56                               | 82.90                        |
| 4          | 44.45   | 0.00   | 99.99                       | 99.56                               | 86.08                        |
| 5          | 44.45   | 0.00   | 99.99                       | 99.56                               | 84.41                        |
| 6          | 44.45   | 0.00   | 99.99                       | 99.56                               | 85.25                        |
| 7          | 44.45   | 0.00   | 99.99                       | 99.56                               | 82.38                        |
| 8          | 44.45   | 0.00   | 99.99                       | 99.56                               | 87.40                        |
| 9          | 44.45   | 0.00   | 99.99                       | 99.56                               | 87.41                        |
| 10         | 52.75   | 0.00   | 99.99                       | 99.55                               | 81.49                        |
| 11         | 55.47   | 0.00   | 99.99                       | 99.57                               | 75.74                        |
| 12         | 44.45   | 0.00   | 99.99                       | 99.56                               | 84.31                        |
| 13         | 50.00   | 0.00   | 99.99                       | 99.56                               | 84.03                        |

## 1.6 Conclusion

This study introduces a postponable acceptance appointment system for a diagnostic clinic. Diagnostic facilities often serve patients of different priority classes. Outpatients are typically scheduled in advance, but higher priority patients (*i.e.*, inpatients and emergency patients) are usually seen as soon as possible. Scheduling of outpatients at the clinic are currently done on an FCFS basis. Thus, high priority outpatients may not receive timely appointments. This challenge motivated us to propose a postponement system in scheduling of different patient classes. The value of the proposed model is that the system can strategically postpone the acceptance of low priority outpatients while waiting for higher priority outpatients. We formulate the problem as a TSSP model in which the first stage estimates the optimal capacity reserved for inpatients and emergency patients. In the second stage, the decisions regarding acceptance and referral of outpatients are made.

Using a data set from the Radiology Department of Prisma Health, we have conducted a series of experiments to test how the model works. The results suggest that postponing the acceptance or referral of outpatient appointment requests up to two days improves the system-wide cost while reducing indirect waiting times. The cost improvement achieved is primarily due to the increase in the utilization of the unused inpatient capacity for outpatients waiting in the queue. In addition, the system prioritizes more urgent outpatients by having them wait only one day in the queue and forcing the less urgent outpatients to wait for two days in the acceptance queue. After analyzing the optimal solutions obtained from our model we developed a simple benchmark policy that can be implemented in real life which performs well.

This study can be extended in multiple directions. For example, in this study we assume

that the duration of visits are constant and identical for each type of patient. Thus, the number of patients that can be seen each day is a fixed number. To consider a more realistic case, uncertain service times can be considered. Furthermore, due to the higher indirect waiting time of lower priority patients, the possibility of no-shows may increase for these classes. Thus, the model can be extended to consider no-shows. Another extension could be to develop a multi-stage stochastic programming approach since the demand uncertainty is revealed over time after each time period. Such multi-stage approaches will be computationally more difficult to solve. Alternatively, the two-stage stochastic program can be used on a rolling horizon basis.

## Chapter 2

# Dynamic Tuberculosis Screening for Healthcare Employees

**Summary:**

### 2.1 Introduction

Tuberculosis (TB) is an infectious disease caused by *Mycobacterium tuberculosis* (Mtb) that mainly affect the lungs, but can also impact other parts of the body. Approximately one third of world's population is infected by Mtb [37]. TB infections are categorized into two groups: latent and active. Individuals with latent infection have the disease but do not show any symptoms, and they are not infectious. However, a latent infection can turn in to active TB. Individuals with active TB are infectious and spread TB. Thus, medication is required for active TB. Although TB infections have been declining in Western Europe and North America, there are groups of people with high prevalence of TB such as immigrants and prisoners [92]. Due to being in contact with these groups, healthcare providers face a higher than usual risk of exposure to TB. Early detection of TB infections is critical to control the spread of the disease. Thus, healthcare workers are usually suggested to take a TB diagnosis test upon initial hiring and thereafter at regular intervals. In fact, CDC guidelines prior to 2005 suggested annual testing for medium or higher risk environments. However, these guidelines were later changed to allow for other policies to be implemented. Yet,

many healthcare facilities such as Prisma Health still follow the old guidelines which usually result in expensive testing for such facilities.

There are currently two common tests for detection of latent and active TB on the market: the skin test and the blood test which is also called Interferon Gamma Release Assay (IGRA). The skin test involves an injection of 0.1 mL of a liquid containing five tuberculin units of purified protein derivative (PPD) into the top layers of skin of the forearm. Once this liquid is injected, the test has to be read within 48-72 hours. A confirmed positive test involves swelling of the injection site, while a negative test has no signs of inflammation. Although the skin test is rapid, it can be prone to false positive results that require expensive further testing. There are different factors that may cause a false positive result such as sensitization of the test to some nontuberculous mycobacteriae, incorrect interpretation of reaction, and incorrect method of the skin test administration [62, 92]. The skin test also shows a false positive if the patient has already had the Bacillus Calmette-Gurin (BCG) vaccine. False positive results cause administration of unnecessary chest X-rays at an extra cost of \$100-\$400 per X-ray [1]. There are also factors such as incorrect interpretation of reaction, some viral illnesses (e.g., measles and chicken pox) and recent live-virus vaccination (e.g., measles and smallpox) that may lead to a false negative skin test result [92]. Individuals who received a false negative skin test may infect others until they begin showing obvious symptoms and begin treatment.

In a TB blood test, a small amount of blood is drawn and sent to a laboratory. Thus, healthcare employees visit the clinic only one time, as opposed to twice for a skin test. The cost of a blood test is higher than that of a skin test (\$10-30 versus \$30-220 [2]), but the test is more accurate. Also, unlike a skin test, the accuracy of a blood test is not affected by prior BCG vaccinations. Although the cost of performing a blood test is much higher than the cost of a skin test, blood tests can still be preferable at least for specific groups of employees when considering the lost time of employees and healthcare professionals who administer the test, as well as the costs related to false positive or false negative results.

In this study, we collaborate with Prisma Health, a healthcare system in South Carolina, that currently requires two skin tests for new employees and an annual skin test for all other employees. Prisma Health does not utilize blood tests currently. As part of our analysis, we categorize the employees into multiple groups based on the risk of infection related to their job, their work environment, birth country and BCG vaccination history. We define an infection rate for each employee

group which depends on the number of infected people the employees in that group potentially get in contact with. We introduce a Markov Decision Process (MDP) model to develop a screening plan for a healthcare facility by determining the type and frequency of TB test to be used for each employee group. The objective is to minimize the expected total cost of the system. Due to the curse of dimensionality, we use Approximate Dynamic Programming (ADP) to find a “near-optimal” solution. Based on this solution, we propose a benchmark policy that is easy to implement by the healthcare facility, and evaluate this simple policy using data obtained from Prisma Health.

## 2.2 Literature Review

There are two streams of literature related to our research. We briefly explain each and highlight our contributions in relation to other studies.

The first stream of studies relates to the analysis of TB screening and its importance for healthcare employees. As mentioned in the previous section, healthcare employees are at a higher than usual risk of getting TB infections. [40] show that even healthcare employees who work in places that are not in direct contact with patients, for instance employees who work in a hospital kitchen, are at a higher risk of being infected. Most healthcare facilities have a regular plan of TB screening for their employees. One common strategy is to administer an annual skin test. [64] tested three different scenarios to propose a TB screening program. The first scenario proposes annual screening for all employees. In the second scenario, only employees with high-risk tasks, such as respiratory therapy, are tested yearly and other employees are tested only after recognized exposure. The third scenario tests all employees only after recognized exposure. They evaluated these scenarios by using both skin and blood tests. They also did a cost effectiveness study on 1000 US healthcare employees with no positive TB history. Results of their experiments indicated that for most US healthcare employees annual TB testing is expensive with limited health gains. Thus, regular annual testing may not be an effective strategy for most health systems.

In addition to the testing frequency, the type of TB diagnosis test to be used is also an important decision. [85] studied the specificity of skin and blood tests among students in a low-tuberculosis incidence setting. They concluded that the blood test performed better in these settings. [48] performed a sensitivity comparison of blood and skin tests in 50 cases with active TB and showed that the sensitivity of the blood test is about 80% while the skin test sensitivity (accuracy ) is about

28%. Accuracy of blood and skin tests in detecting latent and active TB is also presented in other studies [57, 35, 56, 50]. However, further cost analysis studies are needed in this area. [25] performed a cost-effectiveness analysis of TB blood and skin tests by considering the direct test costs and cost of missed work time. [91] extended this study and considered the performance of each test in calculation of total cost of taking a blood test versus a skin test. Thus, cost of subsequent tests and treatments in the case of getting a positive result was included in the total cost. Neither of these studies considered the potential costs that a false negative result may cause. Previous studies also did not analyze the cost-effectiveness of the tests with respect to employee characteristics.

The second stream of relevant studies is on the application of MDP models in the context of prevention, screening, and treatment of diseases. These decisions are typically made sequentially over long periods in uncertain environments [67]. In addition to the patient’s current health status, the uncertainty in progression of the disease, impact of the treatment on the patient, and accuracy of the test results have to be considered in determining the treatment decision [82]. Using MDPs is often appropriate to analyze such problems since the decisions are made sequentially over time in a fundamentally stochastic environment. [30] developed a MDP to model adverse drug reactions in medication treatment of type 2 diabetes. MDPs are also used in breast cancer screening [59, 21, 14, 20], treatment of HIV [77], and public policy decisions related to the transmission of communicable diseases [46, 93]. [14] used a partially observable Markov decision process (POMDP) to take individualized mammography screening decisions while some personal risk features in addition to age and screening history of each patient is considered. They show that by considering these strategies the number of false positive results decreased and quality-adjusted-life-years (QALYs) are improved. This model is extended by [20] and resource constraints are added. They show that allocating capacity efficiently among individuals with different cancer risk levels leads to significant QALYs gains.

Another important application of MDPs in treatment of the diseases is liver and kidney transplant decisions [11, 10, 76, 12]. [11] created an infinite horizon MDP to determine when a patient with end-stage liver disease such as hepatitis C accepts a living-donor transplant. Depending on the quality of the match with the donor and the current health status of the patient, the model determines whether the transplant increases the expected total lifetime of the patient and whether the transplant should be done.

Our contributions to the literature are as follows: 1) We offer a mathematical model in the



area of TB screening, a first in the relevant literature. All previous studies apply simulation or cost evaluation methods that compare the cost of different screening methods. In other words, there is no optimization model to determine TB screening plans. 2) To the best of our knowledge, our study is the first that categorize healthcare employees based on factors that affect the results of TB screening and finds the best test for each group to minimize the expected total cost by considering the infection rate of healthcare employees. 3) Our study is the first to develop a discrete time, infinite horizon MDP model to determine the optimal time between tests for each healthcare employee group in addition to detecting the best test type for each group. 4) For the first time in the literature, we consider the potential impacts and costs of a false negative result (i.e., probable spread of disease) in our formulation.

## 2.3 An MDP Model for the TB Test Scheduling Problem

We propose a discrete time infinite horizon MDP model to formulate the problem. The decision epochs, the state space, the action set, transition probabilities, and cost parameters are described below.

### 2.3.1 Decision Epochs

The decisions on whether or not an employee group should take a TB test and if so what type of test they should take are made annually. Each year, new employees who started their work in the hospital have to take either the blood or the skin test. Current employees with no TB history are also eligible to take a TB test.

### 2.3.2 The State Space

The employees are classified based on their salaries and risk of infection. Considering risk for classification purposes is perhaps obvious, but considering salary groups may not be. Unlike the other studies in the literature, our model captures the opportunity cost of lost time by employees in deciding which test to administer and how often. Thus, considering salary is important in capturing this opportunity cost. Let  $\mathcal{I} = \{1, \dots, I\}$  be the set of employee types based on salary, and  $\mathcal{J} = \{1, \dots, J\}$  be the set of employee types based on infection risk. The infection risk groups

are categorized based on the employees' work locations, the specific job they do, and their BCG vaccination history.

The state of the system is determined by the number of new and ongoing employees with no positive TB history, because the result of the test for employees with positive TB history will almost always be positive. In Prisma Health, these employees with positive TB history fill up a questionnaire each year, and a decision regarding the necessity of an X-Ray test is made based on their responses. We let  $y_{ij}^t$  be the number of current (not new) employees of salary group  $i \in \mathcal{I}$  and risk group  $j \in \mathcal{J}$  who are still employed at the healthcare facility during year  $t$ . The new employees in salary group  $i$  and risk group  $j$  who join the health systems in year  $t$  are represented by  $x_{ij}^t$ .

Employees in different risk groups have different rates of TB infection. The rate of infection of each group in each year depends on the undetected infected employees in each of the groups in the previous year. In particular, undetected infected employees spread TB among other employees and increase the infection rate. Thus, we need to keep track of the number of undetected infected employees. Thus, let  $u_{ij}^t$  be the number of undetected infected employees of salary group  $i$  and risk group  $j$  in the beginning of year  $t$ . Thus, the state space takes the form

$$\vec{s}^t = (\vec{x}, \vec{y}, \vec{u}) = (x_{ij}^t; y_{ij}^t; u_{ij}^t), \quad i \in \mathcal{I}; j \in \mathcal{J} \quad (2.1)$$

We let  $M_{ij}^x$ ,  $M_{ij}^y$  and  $M_{ij}^u$  be the maximum value for current number of employees, new arrivals, and infected employees of salary group  $i$  and risk group  $j$ , respectively. Thus, The state space is finite since  $x_{ij}^t$ ,  $y_{ij}^t$ , and  $u_{ij}^t$  are bounded.

### 2.3.3 The Action Set

The actions that are taken in each state determine whether to administer a test and type of the test to be administered for each employee group. Recall that new employees have to take a test in their first year of employment. Let  $a_{sij}^{xt}$  and  $a_{bij}^{xt}$  denote the decisions regarding the type of the test for new employees of salary group  $i$  and risk group  $j$  in year  $t$  where  $a_{sij}^{xt} = 1$  if skin test is selected and 0 otherwise. Similarly,  $a_{bij}^{xt} = 1$  if blood test is selected and 0 otherwise. Clearly, one of the tests must be selected for new employees. Thus,

$$a_{sij}^{xt} + a_{bij}^{xt} = 1, \quad i \in \mathcal{I}; j \in \mathcal{J}. \quad (2.2)$$

Let  $a_{sij}^{yt}$  and  $a_{bij}^{yt}$  be the decisions regarding the type of the test for current employees of salary group  $i$  and risk group  $j$  at time  $t$ , where  $a_{sij}^{yt} = 1$  if skin test is selected and 0 otherwise. Similarly,  $a_{bij}^{yt} = 1$  if blood test is selected and 0 otherwise. Obviously, if both  $a_{sij}^{yt}$  and  $a_{bij}^{yt}$  are zero, the employees of salary group  $i$  and risk group  $j$  do not take the test at time  $t$ . Thus,

$$0 \leq a_{sij}^{yt} + a_{bij}^{yt} \leq 1, \quad i \in \mathcal{I}; j \in \mathcal{J}. \quad (2.3)$$

The action set takes the form

$$\vec{a}_s^t = (a_{sij}^{xt}, a_{bij}^{xt}, a_{sij}^{yt}, a_{bij}^{yt}), \quad i \in \mathcal{I}; j \in \mathcal{J}, \quad (2.4)$$

and the set of allowable actions in each state  $\vec{s} \in S$  that satisfy constraints (2.2) and (2.3) are denoted by  $A_{\vec{s}}$ .

### 2.3.4 Transitions

Since a portion of employees stop working at the healthcare system during each year, we define  $l_{ij}^t$  as the number of employees of salary group  $i$  and risk group  $j$  who leave the healthcare facility in year  $t$ , which is defined as  $l_{ij}^t \sim \text{Binomial}(y_{ij}^t, p_{ij}^l)$ , where  $p_{ij}^l$  is the probability that an employee of salary group  $i$  and risk group  $j$  leaves the system. Thus,  $l_{ij}^t$  represents the number of employees with no positive TB history who leave.

The new employees who are selected to take skin test, have to take the test twice. The second test must be taken within 7-21 days of the first test. Let  $d_{sij}^{xt}$  and  $d_{sij}^{yt}$  be the total number of new and ongoing employees of salary group  $i$  and risk group  $j$  who take the skin test in year  $t$ . Since there is no difference between the testing protocols for new and ongoing employees if they are selected to take a blood test, we let  $d_{bij}^t$  be the total number of employees of salary group  $i$  and risk group  $j$  who take the blood test in year  $t$ . We define

$$d_{sij}^{xt} = a_{sij}^{xt} x_{ij}^t, \quad i = 1, \dots, I; j = 1, \dots, J, \quad (2.5)$$

$$d_{sij}^{yt} = a_{sij}^{yt} (y_{ij}^t - l_{ij}^t), \quad i = 1, \dots, I; j = 1, \dots, J, \quad (2.6)$$

$$d_{bij}^t = a_{bij}^{xt} x_{ij}^t + a_{bij}^{yt} (y_{ij}^t - l_{ij}^t) \quad i = 1, \dots, I; j = 1, \dots, J. \quad (2.7)$$

Each year, a random number new infections occur in each employee group. Let  $\alpha_{ij}^t$  be the infection probability for an employee of salary group  $i$  and risk group  $j$  in year  $t$ . The infection probability of each group in a year depends on the proportion of undetected infected employees in all employee groups in the previous year, the likelihood that the employees of the group contact with employees of different groups, and the transmission probability conditional on such contact. The infection probability also depends on the percentage of TB infected patients who visit the healthcare facility, probability that an employee of this group contacts an a patient and the transmission probability conditional on such contact. Let  $u_{ij}^t$  denote the total number of undetected infected employees of salary group  $i$  and risk group  $j$  in year  $t$ . We define  $\rho_{ij,i'j'}$  as the probability that an employee of salary group  $i$  and risk group  $j$  contacts individuals of salary group  $i'$  and risk group  $j'$  during that year, and  $\xi_{ij}$  as the TB transmission probability of an employee of salary group  $i$  and risk group  $j$  in case of contacting an infected individual. We define  $\beta$  as the proportion of infected patients who visit the healthcare facility,  $\nu_{ij}$  as the probability that an employee of salary group  $i$  and risk group  $j$  contacts a patient. Thus, we define  $\alpha_{ij}^t$  as

$$\alpha_{ij}^t = \sum_{i'=1}^I \sum_{j'=1}^J \rho_{ij,i'j'} \xi_{ij} \frac{u_{i'j'}^{t-1}}{x_{i'j'}^{t-1} + y_{i'j'}^{t-1}} + \beta \nu_{ij} \xi_{ij} \quad i = 1, \dots, I; j = 1, \dots, J \quad (2.8)$$

We define  $n_{sij}^{xt}$  and  $n_{sij}^{yt}$  as the number of infected new and ongoing employees of salary group  $i$  and risk group  $j$  who have to take the skin test in year  $t$ . We also define  $n_{bij}^t$  as the number of infected employees of salary group  $i$  and risk group  $j$  who take the blood test at time  $t$ . We let  $n_{sij}^{xt}$ ,  $n_{sij}^{yt}$  and  $n_{bij}^t$  follow binomial distributions in the forms of  $Binomial(d_{sij}^{xt}, \alpha_{ij}^t)$ ,  $Binomial(d_{sij}^{yt}, \alpha_{ij}^t)$  and  $Binomial(d_{bij}^t, \alpha_{ij}^t)$ , respectively.

Due to possible false negative test results, a portion of infected employees are undetected. The probability of getting a false negative result depends on the BCG vaccination history of employees. Since BCG vaccination is one of the factors that we consider in defining risk groups, the probability of getting a false negative result depends on the employees' risk groups. Let  $p_{sj}^n$  and  $p_{bj}^n$  be the probability of getting a false negative result for an employee of risk group  $j$  in skin test and blood test, respectively. The number of false negative results in skin test for ongoing employees of salary group  $i$  and risk group  $j$  at time  $t$ ,  $u_{sij}^{yt}$ , follows binomial distribution in the form of  $Binomial(n_{sij}^{yt}, p_{sj}^n)$ . For new employees who get tested with a skin test, getting a false negative result means getting false negative results in both the initial and the follow up tests. Thus, the number

of false negative results  $u_{sij}^{xt}$ , follows a distribution in the form of  $Binomial(Binomial(n_{sij}^{xt}, p_{sj}^n), p_{sj}^n)$ . We let  $u_{bij}^t$  be the number of false negative results for employees of salary group  $i$  and risk group  $j$  in year  $t$  who have taken blood test which follows binomial distribution  $Binomial(n_{bij}^t, p_{bj}^n)$ .

The number of new and ongoing employees of salary group  $i$  and risk group  $j$  who have got true positive skin test results are defined as  $q_{sij}^{xt}$  and  $q_{sij}^{yt}$ , respectively, and calculated as  $q_{sij}^{xt} = n_{sij}^{xt} - u_{sij}^{xt}$  and  $q_{sij}^{yt} = n_{sij}^{yt} - u_{sij}^{yt}$ . Similarly,  $q_{bij}^t$  represents the total number of employees of salary group  $i$  and risk group  $j$  who have true positive blood test results in year  $t$  where  $q_{bij}^t = n_{bij}^t - u_{bij}^t$ .

According to our definitions, sum of  $u_{sij}^{xt}$ ,  $u_{sij}^{yt}$  and  $u_{bij}^t$  shows the number of undetected infected employees in salary group  $i$  and risk group  $j$  that were given one of the TB tests in year  $t$ . There might also be infected employees in the employee groups that were not given any of the tests in year  $t$ . Let  $d_{nij}^t$  be the number of employees in salary group  $i$  and risk group  $j$  that was not tested in year  $t$ , which is defined as

$$d_{nij}^t = (1 - a_{sij}^{yt} - a_{bij}^{yt})(y_{ij}^t - l_{ij}^t), \quad i = 1, \dots, I; j = 1, \dots, J. \quad (2.9)$$

The number of infected employees in these groups of employees also follows binomial distribution with infection probability of each group. Let  $u_{nij}^t$  be the number of infected employees of salary group  $i$  and risk group  $j$  that was not tested in year  $t$  where,  $u_{nij}^t \sim Binomial(d_{nij}^t, \alpha_{ij})$ . We define  $u_{ij}^t$  as the total number of undetected infected employees where

$$u_{ij}^t = u_{sij}^{xt} + u_{sij}^{yt} + u_{bij}^t + u_{nij}^t, \quad i = 1, \dots, I; j = 1, \dots, J. \quad (2.10)$$

Some uninfected employees might receive false positive results from the TB tests, which leads to a follow up X-ray at extra cost. We let  $r_{sij}^{xt}$  and  $r_{sij}^{yt}$  be the number of new and ongoing employees of salary group  $i$  and risk group  $j$  who have received false positive skin test results. Similarly, we let  $r_{bij}^t$  show the number of employees of salary group  $i$  and risk group  $j$  who have got false positive blood test results. These variables follow the distributions  $r_{sij}^{xt} \sim Binomial(d_{sij}^{xt} - n_{sij}^{xt}, p_{sj}^p) + Binomial(d_{sij}^{xt} - n_{sij}^{xt} - Binomial(d_{sij}^{xt} - n_{sij}^{xt}, p_{sj}^p), p_{sj}^p)$ ,  $r_{sij}^{yt} \sim Binomial(d_{sij}^{yt} - n_{sij}^{yt}, p_{sj}^p)$  and  $r_{bij}^t \sim Binomial(d_{bij}^t - n_{bij}^t, p_{bj}^p)$  where  $p_{sj}^p$  and  $p_{bj}^p$  are the probabilities that an employee of salary group  $i$  and risk group  $j$  gets a false positive result in skin and blood tests, respectively.

Once the testing decisions are made, the stochastic elements that determine the state tran-

sition are the new arrivals of employees, the employees who left the system, and the total number of infections. The evolution of state space elements are captured by the following equations:

$$y_{ij}^{t+1} = y_{ij}^t + x_{ij}^t - l_{ij}^t - n_{sij}^{xt} - n_{sij}^{yt} - n_{bij}^t - u_{nij}^t, \quad i = 1, \dots, I; j = 1, \dots, J \quad (2.11)$$

$$u_{ij}^{t+1} = u_{sij}^{xt} + u_{sij}^{yt} + u_{bij}^t + u_{nij}^t, \quad i = 1, \dots, I; j = 1, \dots, J \quad (2.12)$$

The trajectory of the system is represented by  $\{(\vec{s}^t, \vec{a}_s^t) : t = 1, 2, \dots\}$ , where  $\vec{s}^t$  is the state of the system and  $\vec{a}_s^t$  is the action that is taken in year  $t$ . The stochastic evolution of the system is represented by  $\vec{s}^{t+1} = F(\vec{s}^t, \vec{a}_s^t, g(\vec{s}^t, \vec{a}_s^t))$ , where  $F(\cdot, \cdot, \cdot)$  is a transfer mapping and  $g(\vec{s}^t, \vec{a}_s^t)$  is a random element that contains all the random quantities in the system at time  $t$ . These definitions are used for defining the value function.

### 2.3.5 The costs

The cost associated with each state-action pair drives from four sources: cost of doing the tests ( $c^b, c^s$ ), cost of doing the X-ray ( $c^x$ ), cost of lost time of employees ( $c_i^l$ ), and cost of undetected infections ( $c_i^u$ ).

$$c(\vec{s}, \vec{a}) = \sum_{i=1}^I \sum_{j=1}^J \left( c^b d_{bij}^t + c^s d_{sij}^t + c^x (q_{sij}^{xt} + q_{sij}^{yt} + q_{bij}^t + r_{sij}^{xt} + r_{sij}^{yt} + r_{bij}^t) + c_{ij}^u u_{ij}^t + c_i^l w_{ij}^t \right), \quad (2.13)$$

where  $w_{ij}^t$  is the expected total time spent at the testing clinic for employees of salary group  $i$  and risk group  $j$  in year  $t$ . For employees who are taking blood test in year  $t$ ,  $w_{ij}^t$  either only depends on the time spent administering the blood test, or also includes the time spent on X-Ray if result of the blood test is positive. However, for employees who are taking the skin test, the total skin test time depends on whether or not the employee should take the test twice or once (i.e., employee is new or not). There is also a possibility that the employee is not able to complete the second step of the skin test within the required time and has to repeat both steps. This also affects the expected total time an employee spends in taking the skin test.

### 2.3.6 Optimality equation

The value function  $v(\vec{s})$  corresponds to the total expected discounted cost for state  $\vec{s}$  over the infinite horizon.

$$v(\vec{s}) = \min_{\vec{a} \in A_{\vec{s}}} \left\{ c(\vec{s}, \vec{a}) + \lambda \mathbb{E}(v(\vec{s}')) \right\}, \quad \forall \vec{s} \in S, \quad (2.14)$$

where the expectation is taken with respect to  $s' = F(\vec{s}^t, \vec{a}, g(\vec{s}^t, \vec{a}))$  and  $\lambda \in [0, 1)$  is a discount factor. To find the optimal policy we need to solve equation (3.20). Since the state space and action set are finite, there exists a stationary optimal policy. However, the size of the state space make a direct solution to (3.20) impractical.

## 2.4 Approximate Dynamic Programming

Due to the curse of dimensionality, we use approximate dynamic programming to estimate the optimal policy. First, we transform the MDP model into its equivalent linear program (LP) as follows:

$$\max \sum_{\vec{s} \in S} \gamma(\vec{s}) v(\vec{s}) \quad (2.15)$$

s.t.

$$c(\vec{s}, \vec{a}) + \lambda \sum_{\vec{s}' \in S} \mathbb{P}(\vec{s}' | \vec{s}, \vec{a}) v(\vec{s}') \geq v(\vec{s}) \quad \forall \vec{s} \in S, \vec{a} \in A_{\vec{s}}, \quad (2.16)$$

where  $\mathbb{P}(\cdot | \cdot, \cdot)$  is the transition probability and  $\gamma(\cdot)$  is the probability distribution over the initial state of the system. The LP formulation does not avoid the curse of dimensionality. Thus, we approximate the value function by using a specific parameterized form where the interactions of employees from different groups are not considered. The resulting approximate value function can be written as the summation of separate value functions for each employee group as follows:

$$v(\vec{s}) \approx \sum_{i=1}^I \sum_{j=1}^J v_{ij}(\vec{s}_{ij}), \quad \forall \vec{s} \in S, \quad (2.17)$$

where  $\vec{s}_{ij}$  represents the state of the system for employees of salary group  $i$  and risk group  $j$  which shows the number of current employees, new arrivals and infected employees of salary group  $i$  and

risk group  $j$ . Moreover,  $v_{ij}(\vec{s}_{ij})$  denotes the value function for employees of salary group  $i$  and risk group  $j$  that is defined as follows:

$$v_{ij}(\vec{s}_{ij}) = \min_{\vec{a}_{ij} \in A_{\vec{s}_{ij}}} \left\{ c^b d_{bij}^t + c^s d_{sij}^t + c^x (q_{sij}^{xt} + q_{sij}^{yt} + q_{bij}^t + r_{sij}^{xt} + r_{sij}^{yt} + r_{bij}^t) + c_i^u u_{ij}^t + c_i^l w_{ij}^t + \lambda \left( \sum_{\vec{s}_{ij}' \in S_{ij}} \mathbb{P}(\vec{s}_{ij}' | \vec{a}_{ij}, \vec{s}_{ij}) \times v_{ij}(\vec{s}_{ij}') \right) \right\}, \quad \forall i = 1, \dots, I, j = 1, \dots, J, \vec{s}_{ij} \in S_{ij}. \quad (2.18)$$

When the LP-based ADP presented here is used, the number of decision variables grows linearly with the number of employees in each group. In contrast, if the original LP formulation is used, the number of variables grows exponentially with the number of employees in each group. Thus, we expect that this approximate reformulation will help addressing the computational challenges due to curse of dimensionality.

Since  $w_{ij}^t$  includes both the service time (time of taking the tests) and waiting time of employees of salary group  $i$  and risk group  $j$  at time  $t$ , states of other groups affect the value of  $w_{ij}^t$ . In the LP-based ADP, we estimate  $w_{ij}^t$  for each group independently of the states of the other employee groups. Thus, we calculate an upper bound for  $w_{ij}^t$  and use this upper bound in our approximation. To compute the upper bound, we assume that all others groups have the largest possible number of employees that we can have in each group, and all employees are taking the skin test. We build a pre-processing simulation model to estimate  $w_{ij}^t$  for each group depending on its current state.

Recall that we have defined  $\gamma(\vec{s}) \forall \vec{s} \in S$  as the probability distribution over the initial state of the system. We let  $\gamma_{ij}(\vec{s}_{ij}) \forall \vec{s}_{ij} \in S_{ij}$  be the marginal initial state distribution for employees of salary group  $i$  and risk group  $j$ , where  $S_{ij}$  is the set of all possible states for each group. Also, let  $\Omega = \{(\vec{s}_{ij}, \vec{a}_{ij}) : \vec{s}_{ij} \in S_{ij}, \vec{a}_{ij} \in A_{\vec{s}_{ij}}, \forall i = 1, \dots, I, j = 1, \dots, J\}$  be the set of all feasible state-action



pairs. The dual of the LP-based ADP model can be expressed as

$$\min \sum_{i=1}^I \sum_{j=1}^J \sum_{(\vec{s}_{ij}, \vec{a}_{ij}) \in \Omega} c(\vec{s}_{ij}, \vec{a}_{ij}) \delta(\vec{s}_{ij}, \vec{a}_{ij}) \quad (2.19)$$

s.t.

$$\sum_{\substack{(\vec{s}'_{ij}, \vec{a}_{ij}) \in \Omega \\ \vec{s}'_{ij} = \vec{s}_{ij}}} \delta(\vec{s}'_{ij}, \vec{a}_{ij}) - \lambda \sum_{(\vec{s}'_{ij}, \vec{a}_{ij}) \in \Omega} \mathbb{P}(\vec{s}_{ij} | \vec{s}'_{ij}, \vec{a}_{ij}) \delta(\vec{s}'_{ij}, \vec{a}_{ij}) \geq \gamma_{ij}(\vec{s}_{ij}), \quad \forall i = 1, \dots, I, j = 1, \dots, J, \vec{s}_{ij} \in S_{ij} \quad (2.20)$$

$$\delta(\vec{s}_{ij}, \vec{a}_{ij}) \geq 0, \quad \forall i = 1, \dots, I, j = 1, \dots, J, (\vec{s}_{ij}, \vec{a}_{ij}) \in \Omega \quad (2.21)$$

In the above formulation,  $c(\vec{s}_{ij}, \vec{a}_{ij})$  shows the immediate cost for employees of salary group  $i$  and risk group  $j$  in state  $\vec{s}_{ij}$  while action  $\vec{a}_{ij}$  is taken. The advantage of solving the dual problem is that the number of constraints is fewer. However, the number of variables is still very large. Thus, we use the column generation algorithm to obtain the optimal solution. The column generation algorithm is started with a small set of feasible state-action pairs (i.e., columns) to the dual problem, which is called the master problem. Then one or more violated constraints in the primal problem are found by solving a subproblem. The state-action pair(s) corresponding to these violated constraints are added to the master problem as new columns. The procedure continues until no primal constraint is violated.

We consider a master problem associated with the formulation (2.19-2.21), and let  $v_{ij}^D(\vec{s}_{ij})$  be the dual variable associated with constraint (2.20) for all  $\vec{s}_{ij} \in S_{ij}$  for  $i = 1, \dots, I, j = 1, \dots, J$ . The corresponding subproblem can be written as:

$$z^{sub} = \min_{(\vec{s}_{ij}, \vec{a}_{ij}) \in \Omega} \left\{ \sum_{i=1}^I \sum_{j=1}^J (c(\vec{s}_{ij}, \vec{a}_{ij}) - v_{ij}^D(\vec{s}_{ij}) + \lambda \sum_{\vec{s}'_{ij} \in S_{ij}} \mathbb{P}_{ij}(\vec{s}'_{ij} | \vec{s}_{ij}, \vec{a}_{ij}) v_{ij}^D(\vec{s}'_{ij})) \right\} \quad (2.22)$$

Formulation (2.22) is a generalized multidimensional knapsack problem. To solve the subproblem, we reformulate it as an integer program that is easy to implement in optimization solvers.

For each employee of salary group  $i$  and risk group  $j$  we let  $\Omega_{ij} = \{(\vec{s}_{ij}, \vec{a}_{ij}) : \vec{s}_{ij} \in S_{ij}, \vec{a}_{ij} \in A_{\vec{s}_{ij}}\}$  be set of all feasible state-action pairs in the master problem and note that  $f_{ij}(\vec{s}_{ij}, \vec{a}_{ij}) = c(\vec{s}_{ij}, \vec{a}_{ij}) - v_{ij}^D(\vec{s}_{ij}) + \lambda \sum_{\vec{s}'_{ij} \in S_{ij}} \mathbb{P}_{ij}(\vec{s}'_{ij} | \vec{s}_{ij}, \vec{a}_{ij}) v_{ij}^D(\vec{s}'_{ij})$  is the corresponding reduced cost for such a state-action pair. Thus, the subproblem finds the feasible state-action pairs such that

$\sum_{i=1}^I \sum_{j=1}^J f_{ij}(\vec{s}_{ij}, \vec{a}_{ij}) < 0$ . We let  $z_{ij}^{\vec{s}_{ij}, \vec{a}_{ij}}$  be a binary variable that is one if action  $\vec{a}_{ij}$  is chosen in state  $\vec{s}_{ij}$  for employees of salary group  $i$  and risk group  $j$ ; and zero, otherwise. Then, we transform the subproblem into a generalized assignment problem which chooses one feasible state-action pair for each employee group. The formulation for the assignment problem is as follows:

$$\min \sum_{i=1}^I \sum_{j=1}^J \sum_{(\vec{s}_{ij}, \vec{a}_{ij}) \in \Omega_{ij}} f_{ij}(\vec{s}_{ij}, \vec{a}_{ij}) z_{ij}^{\vec{s}_{ij}, \vec{a}_{ij}} \quad (2.23)$$

subject to

$$\sum_{(\vec{s}_{ij}, \vec{a}_{ij}) \in \Omega_{ij}} z_{ij}^{\vec{s}_{ij}, \vec{a}_{ij}} = 1, \quad i = 1, \dots, I, j = 1, \dots, J \quad (2.24)$$

$$z_{ij}^{\vec{s}_{ij}, \vec{a}_{ij}} \in \{0, 1\}, \quad i = 1, \dots, I, j = 1, \dots, J \quad \forall (\vec{s}_{ij}, \vec{a}_{ij}) \in \Omega_{ij} \quad (2.25)$$

The objective function (2.23) chooses the column with minimum reduced cost. This new formulation can be implemented and solved using commercial optimization solvers.

## 2.5 Numerical Study and Results

In this section, we implement the approximate dynamic programming scheme described in the previous section using data obtained from one of the Prisma Health hospitals and we obtain the optimal TB screening policy. Based on the optimal policy we propose some benchmark policies and simulate them to show how these policies improve the system compare to the current policies.

### 2.5.1 Input data

In this study we consider three employee salary groups; (i) physicians, (ii) nurses and (iii) other employees, and three employee risk groups; (i) BCG vaccinated employees, (ii) employees who work in high risk locations, (iii) and employees who work in low risk locations. Thus, a total of nine groups are presumed. We need to mention that the yearly arrivals of employees are assumed to be stochastic and modeled as truncated Poisson processes. Yearly arrival rates, leaving probabilities, probabilities of contacting patients and employees, TB transmission probabilities, TB infection probability, false positive and false negative probabilities, and cost parameters are the parameters that are used in the MDP model. The estimation for these parameters are presented Table 3.

Probabilities of contacting patients for different employee groups  $\nu_{ij}$ , are estimated based on our discussions with Prisma Health physicians. In using ADP, we assume that employees of different groups do not contact with each other and employees who are in the same group definitely contact each other. Thus,  $\rho_{ij,i'j'}$ , which is defined as the probability that an employee of salary group  $i$  and risk group  $j$  contacts individuals of salary group  $i'$  and risk group  $j'$ , is zero if  $i \neq i'$  or  $j \neq j'$  and it is one for employees in the same group. We base our estimates of TB Transmission probabilities,  $\xi_{ij}$ , for each group, we on the study by [5]. Since Prisma health did not have exact information about the percentage of patients with TB who come to the hospital,  $\beta$ , we used the information from [64] to estimate this parameter. In that study, the average percentage of individuals infected with TB among healthcare workers is reported to be about 2%. Thus, we simulated and calibrated the system under the hospital's current policy (i.e., annual skin tests for all employees) so that the value of  $\alpha_{ij}^t$  in the long run is 2%. The results of this simulation showed that value of  $\beta$  which brings %2 for  $\alpha_{ij}^t$  is 0.1. Hence, we use this value in our experiments. The false positive and false negative probabilities of the tests are taken from [29, 50, 28, 42]. The remaining parameters used in the model (i.e., the cost parameters) are estimated based on the hospital data.

### 2.5.2 Experimental setup

The optimization problem is implemented in C++. The subproblems and master problems are solved on an Intel Core i7-9700 CPU utilizing the Gurobi 9.0 solver. The computational time of solving the total problem for one instance is about 120 minutes.

### 2.5.3 Optimal policy for the base model

The optimal policy which dictates testing actions in each state is obtained by solving the ADP model. After analyzing the optimal policy, we observed that changes in number of new arrivals does not have huge impact on the optimal testing decision. Thus, to visualize the optimal policy and make its discussion easier, we created plots of current number of employees ( $y_{ij}^t$ ) versus number of infected employees ( $u_{ij}^t$ ). These plots are presented as Figures 2.1, 2.2, and 2.3.

Each plot has regions that are numbered (1), (2) or (3). Region (1) shows the states where the optimal action is not taking any test. Regions (2) and (3) represent the states where the optimal action is taking skin test and blood test, respectively. In the following we briefly describe these

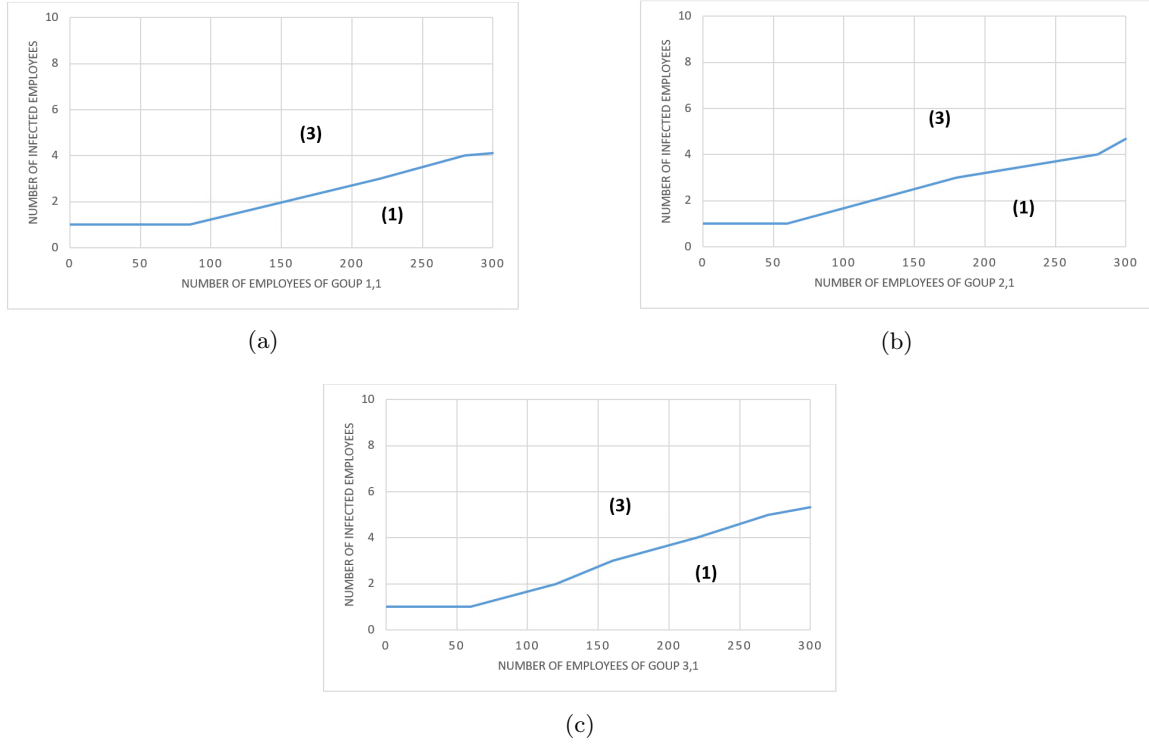


Figure 2.1: Optimal policy for risk group 1

graphs.

Figure 2.1 shows the optimal policies for employees of risk group 1 (i.e., employees who had the BCG vaccination). We observe that for employees across all salary groups who had the BCG vaccination, the optimal decision is either no testing or the blood test depending on the level of infection spread among the employees. This is because the skin tests are less cost-effective for these employees due to the major false positive risk.

Figure 2.2 depicts the optimal policies for employees of risk group 2 (i.e., employees who work at low risk locations). The results show that for physicians, the optimal decision is either no test or blood test. That is due to the high cost of lost time incurred when these employees get skin tests. On the other hand, for nurses and other employees depending on the number of employees and number of infected ones the decision would be either no test, skin test or blood test.

Finally, the optimal policies for risk group 3 (i.e., employees who work at high risk locations) is shown in Figure 2.3. As we expected, for physicians the optimal decision is either no test or blood test, and for nurses and others the optimal decision can be no test, skin test or blood test depending

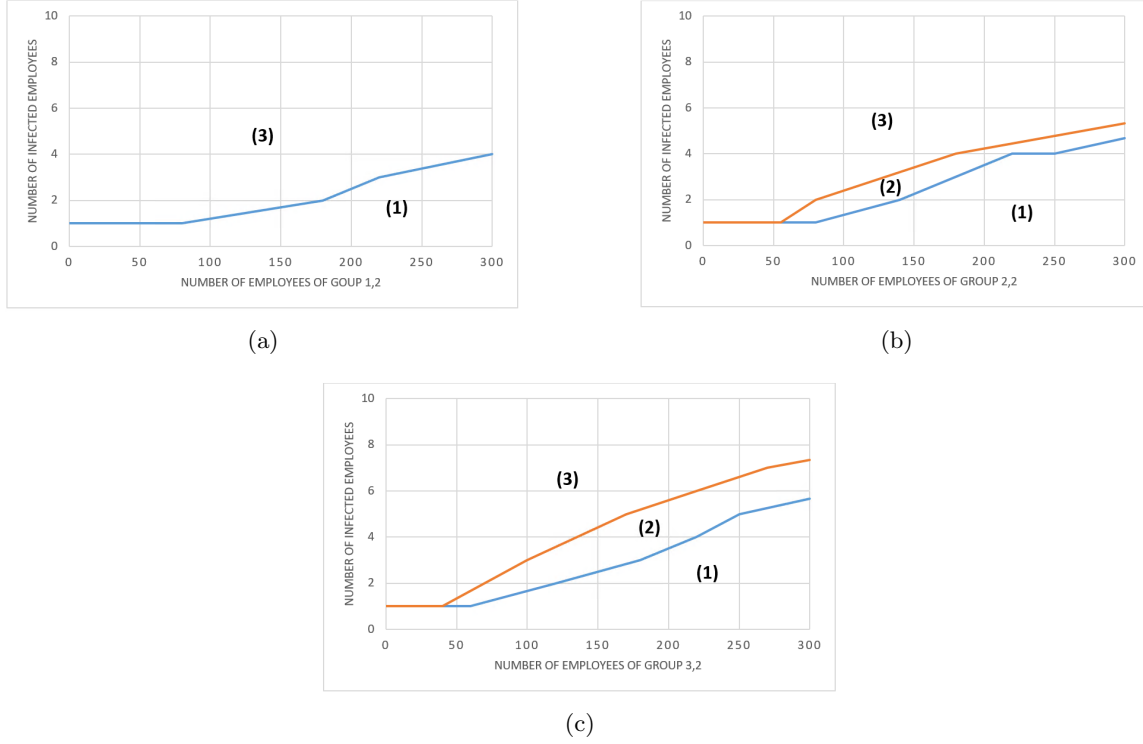


Figure 2.2: Optimal policy for risk group 2

on the level of infection spread.

#### 2.5.4 Simulation of the optimal policy

The Optimal policy indicates the optimal action (no test, skin test, or blood test) for each employee group in different states. However, how often the tests are performed under the optimal policy is not readily known without further analysis. we are interested in knowing these frequencies since the current practice as well as CDC guidelines are based on testing frequency and not on the number of employees and infections detected. Thus, we simulate the optimal policy to estimate the time between tests for each employee group using the number of times that each group visits the states that require testing. The simulation model is implemented in C++ and ran for 100 years. The results of the simulation are presented in Table 2.1. We report the optimal action which is the type of the test administered and frequency of that type of test for each employee group. As mentioned above, the frequency of the tests is estimated based on the number of times that the test is administrated in a 100 year horizon..

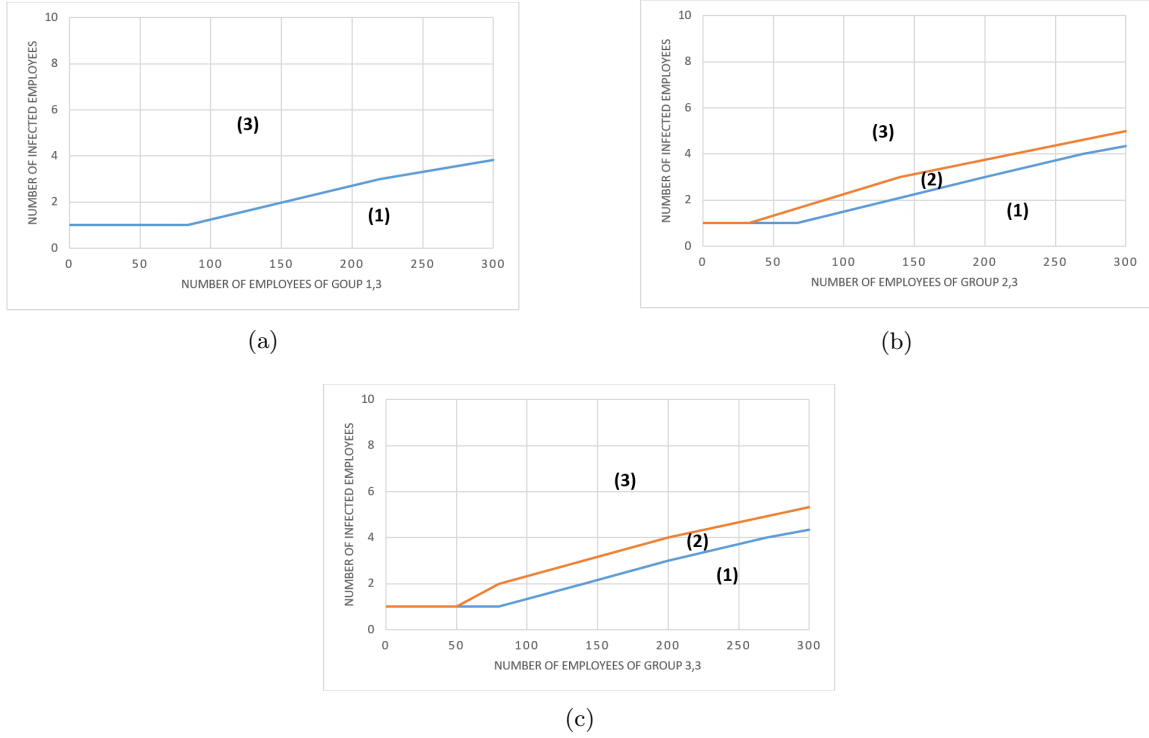


Figure 2.3: Optimal policy for risk group 3

For some employee groups, the frequency of testing was very low, since the groups visited states that require testing infrequently. Thus, instead of just providing a recommendation of practically never testing these employees, we define the critical ratio between number of infected employees and current number of employees in these groups that would trigger testing. In other words, we define a threshold for percentage of infected employees where for larger percentages either a blood test or a skin test is required. These thresholds are extracted from the optimal policy.

In order to observe the advantages of the proposed policy given in Table 2.1, we simulate this policy and compare the result with the hospital's current policy which is administering skin tests every year for all employee groups. These simulations are also implemented in C++. In these simulations, we still assume that employees of different groups do not interact. Table 2.2 shows the corresponding results. The proposed policy decreases the average screening cost of the healthcare facility by optimizing the type and frequency of the tests for different employee groups while the average infection rate does not increase. Based on these results, it is beneficial for the hospital to update their screening guidelines for different employee groups instead of annually testing all

Table 2.1: Optimal actions based on the simulation

| Employee group | Optimal action for new employees | Optimal action for current employees  |
|----------------|----------------------------------|---|
| 1,1            | take blood test                  | take blood test every year  |
| 1,2            | take blood test                  | take blood test every year  |
| 1,3            | take blood test                  | take blood test every year  |
| 2,1            | take blood test                  | infrequently - take blood test if percentage of infected employees is greater than 1.6%   |
| 2,2            | take blood test                  | take blood test every 3 years   |
| 2,3            | take blood test                  | take blood test every 3 years   |
| 3,1            | take blood test                  | infrequently - take blood test if percentage of infected employees is greater than 1.7%   |
| 3,2            | take blood test                  | infrequently - take blood test if percentage of infected employees is greater than 2.2% and skin test if it is between 1.7 and 2.2% |
| 3,3            | take blood test                  | take blood test every 2 years   |

Table 2.2: Comparison of the proposed policies and the current policy

|                        | Current policy | Proposed policy | Optimal policy |
|------------------------|----------------|-----------------|----------------|
| Average yearly cost    | 17417 \$       | 8752 \$         | 6750 \$        |
| Average infection rate | 1.5 %          | 1.5 %           | 1.2 %          |

employees regardless of risk and salary characteristics.

## 2.6 Conclusion

This study includes a dynamic TB screening model for healthcare employees. Although, rate of TB infection in the US is low, healthcare employees are still in risk of getting infected due to being in direct contact with patients. Thus, regular screening seems to be an effective preventive approach. For screening purposes, hospitals can administer skin or blood tests. These tests bring costs for the healthcare facility. Thus, providing an optimal screening policy that minimizes the infection rate and cost of the healthcare facility is important.

This challenge motivated us to propose a dynamic model that specifies the best TB test and its frequency for each group of employee. We formulated this problem as a MDP and used ADP to estimate the solution. We used our formulation and a data set obtained from one of the Prisma Heealth hospitals to determine the optimal TB screening test for each employee group. Currently, an annual skin test is required in this hospital. Furthermore, we simulated the optimal policy to estimate the optimal screening frequency for each group of employees. Then, we simulated these resulting policy and compare the results with the current screening policy implemented at our partnering hospital. Results confirm the improvement in the average cost of the healthcare facility without an increase in the observed infection rate.

This study can be extended in multiple directions. Our model does not differentiate between the latent and active TB infections. However, latent and active TB may have different impacts on the results of the tests or may affect the spread of the disease differently. Furthermore, by using ADP we assumed that employees of different groups do not contact with each other. However, to consider a more realistic case, contacting with different groups can be considered.



## Chapter 3

# Risk Based Staffing for Pandemic Response

**Summary:** Protection of healthcare employees is of paramount importance in controlling the spread of disease during a pandemic such as the one caused by Coronavirus Disease 2019 (COVID-19). While millions of people around the world stay home to minimize the spread of COVID-19, healthcare employees do the opposite as they work longer shifts and attend to more patients than normal. These work conditions increase the risk of infection for them. Thus, healthcare systems are looking for strategies to protect their employees in settings where social distancing cannot be maintained. In this study, we propose a staffing strategy where the healthcare workforce is organized into teams and the teams are scheduled so that only some of them work in each time period to decrease the risk of infection spread. We create a discounted Markov decision process (MDP) model to decide how teams should be assigned to work when maximizing the number of working teams and minimizing infection rate of employees is desired. Due to the large size of state space of the MDP, we use state reduction techniques and the policy iteration algorithm to solve the problem. The results show that problem size is greatly reduced by our state space reduction algorithm and the optimal policy obtained for the approximated problem can improve the objective function value considerably compared to a benchmark policy that has all teams working all the time.

### 3.1 Introduction

COVID-19 is an ongoing pandemic, heavily affecting almost all nations throughout the globe. The virus spreads so rapidly that two weeks after the first cases were diagnosed 1000 patients tested positive globally [81]. With increasing infection and hospitalization rates, the pandemic is putting a significant load on healthcare systems. Hospitals need to maintain and even increase their capacity to be able to effectively handle the high workload created by the pandemic. Furthermore, healthcare employees are at a higher risk of being infected due to their interactions with symptomatic and asymptomatic COVID-19 patients and not being able to adhere to standard social distancing measures. Figures from China’s National Health Commission show that more than 3300 healthcare employees have been infected as of early March 2020 [53]. In Italy, 20% of healthcare employees who had contacted COVID-19 patients were infected, and some have died [53].

Spread of infection among the physicians can have a serious detrimental effect on pandemic response efforts due to reduced workable-physician-days. Thus, it is of critical importance to make staffing and scheduling decisions in a way that would limit exposure of healthcare employees to infection (and hence, minimize their risk) while the care delivered to the patients is still guaranteed to be timely and of high quality.

There are four main ways infectious diseases, and in particular, COVID-19, can spread among healthcare employees in a hospital setting: (i) through interactions with symptomatic patients, (ii) through interactions with asymptomatic patients, (iii) through interactions among healthcare employees, and (iv) through interactions with infected surfaces. Physicians can use personal protective equipment (PPE) while caring for symptomatic patients to reduce the risk of infection [34]. However, infection prevention through consistent use of PPE can be difficult when interacting with asymptomatic patients and other healthcare employees, or using common areas and equipment. Thus, additional operational measures on how the employees work and interact are needed to eliminate or control these possible infection routes in hospitals.

Recent studies show that grouping healthcare employees into teams and scheduling them to work in “bubbles” instead of independently without limiting their interactions decreases the infection rate among healthcare employees during the COVID-19 pandemic [75, 71]. In this study, we build a dynamic model for organizing and employing the employees as teams. We introduce a Markov Decision Process (MDP) model to find the optimal policy for weekly scheduling of these teams. In

our formulation, besides the infection probability of the general population, we introduce an infection probability for each employee team. This probability may increase due to teams being scheduled to work during each week. Our objective is to maximize the number of working employees at each week by scheduling at least the minimal required capacity and controlling the infection probability in employee teams at the same time. Due to the large size of the problem, we use a state reduction technique to decrease the size of the state space. Then, we use the policy iteration algorithm to find the optimal policy for the new approximated system. We solve our model with publicly available COVID-19 infection data from South Carolina Department of Health and Environmental Control and operational data that is given by our healthcare system of interest, Prisma Health in South Carolina.

## 3.2 Literature Review

There are three streams of literature related to our study. In the following, we describe each and mention our contribution with respect to other studies. The first stream of studies relates to the infection control for healthcare employees and organization of workforce during pandemics. Healthcare employees play an important role in controlling the spread of infectious diseases. Moreover, they are at high risk of infection because of the nature of their work. In case of availability, vaccination is the first choice of infection control for healthcare employees in pandemics [22, 13, 83]. Besides vaccination, studies found that personal protective equipment and antiviral treatments are helpful [83, 84]. [24] provide guidance to provide practicable help for the control of parvovirus B19 infection in healthcare settings. This guidance focuses on contact of at-risk employees (such as pregnant women) with patients and employees with at-risk patients.

Besides controlling infection via vaccinations and the use of PPE, organizing the workforce to mitigate infection risk is another critical task during pandemics. Agent-based simulation models are commonly used in many epidemiological studies on infectious diseases [23, 49, 41]. [23] used an agent-based model to simulate the effect of influenza pandemic on healthcare settings and employees. Their model forecasts how many employees may become infected, how many employees will be available, how many patients each employee should see, and how soon the employees should be protected. [75] simulate multiple strategies to organize healthcare workforce during pandemics and present an application for the COVID-19 pandemic. They introduce the scenario of desynchronization of the

workforce in which teams are dichotomized and each half of the team works on alternating weeks. They simulate this scenario and compare the results with the case where there is just one team of employees. Based on the results, having two employee teams decreases the infection rate in the healthcare setting. In this study, we expand the idea of having employee teams. However, our study and contributions differ from [75] study in several aspects. First, our model allows for more than two teams. Second, In [75] study, a simulation model is used to compare one team versus two teams of employees, and there is a fix strategy for scheduling them (i.e., one team is working, the other team is resting). In contrast, we build an optimization model to identify the optimal scheduling policy for employee groups and determine which groups have to work depending on the state of the system. Third, our models account for a minimum of two-weeks quarantine for infected employees which is a more accurate infectious period estimate for COVID-19 than the one week considered in [75]. Fourth, using the spread pattern of COVID-19 that is characterized based on the real data, we determine the infection rates for healthcare employees as a function of the infection rate of the general population and the interactions of the employees with patients and with each other.

The second stream of studies related to ours is application of MDP models in the context of prevention, screening, and treatment of diseases. These decisions are often taken in an uncertain environment [67]. Moreover, to make a decision regarding treatment of the disease, factors other than the patient’s current health status may need to be considered. For example, the uncertainty in progression of the disease and impact of the treatment on the patient [82] can be considered. MDP models are appropriate tools in characterizing these problems as they provide ways to model such features. For instance, [30] employed a MDP to model adverse drug reactions in some medications that are used to treat patients with type 2 diabetes. MDPs are also used in breast cancer screening [59, 21, 14, 20], cervical cancer screening [8, 65] and the treatment planning for HIV [77]. [14] developed a partially observable Markov decision process (POMDP) to take individualized mammography screening decisions while some personal risk features in addition to age and screening history of each patient is considered. They show that by considering these strategies the number of false positive results decreased and quality-adjusted-life-years (QALYs) are improved. [20] extended their work and added resource constraints. They show that allocating capacity efficiently between individuals with different cancer risk levels leads to significant QALYs gains. MDPs are also used in problems with liver and kidney transplant decisions [11, 10, 76, 12]. [12] formulated their problem as a discrete-time, infinite-horizon Markov decision process model. Their model determines whether the patient

should have a transplantation now or wait. Additionally, in case of a transplantation decision, the model also addresses whether for a patient with end-stage liver disease a cadaveric liver from a deceased donor or a portion of a living-donor’s liver has to be transplanted.

The third stream of studies focus on disease characteristics and mechanisms of spread for COVID-19. We use this emerging literature of related studies to inform our model structure and parameters. Spread of COVID-19 is a crisis all over the world. It is shown that the transmission of the disease happens by coughing and sneezing by symptomatic patients and also from asymptomatic individuals before symptoms are observable [74]. Patients may also be infectious on clinical recovery [79]. The COVID-19 virus can also stay on surfaces for couple of days and infect individuals [44]. According to the [55] study, the incubation period of the virus is estimated as 5.1 days, and individuals who develop the symptoms usually will do so within 11.5 days. This results were obtained by studying 181 confirmed cases in Hubei province, China. The fact that healthcare employees are in higher risk of getting infected is shown in [89, 53, 66]. Using PPE, provision of food, having isolation room for them and rest are some operations that may help to decrease this risk for healthcare employees [53, 34]. However, more strategic actions may be required to control infection rate for healthcare employees.

In this study, we investigate the benefits of controlled segregation and isolation of healthcare employees in controlling COVID-19 spread in healthcare facilities. We considered teams of physicians that work in the same specialty. To the best of our knowledge, we are first to propose a discrete time infinite horizon discounted MDP model to organize the healthcare workforce in a pandemic. Furthermore, ours is the first study to develop an optimization model for infection control among healthcare employees through strategic scheduling practices. We propose policies for scheduling of teams of healthcare employees that maintains a sufficient level of utilization of the workforce while mitigating the infection risk.

### 3.3 An MDP Model for Staff Scheduling

The preliminary results obtained using ABM indicate that scheduling teams of healthcare employees with limited mixing can be beneficial in decreasing the risk of COVID-19 infections among the said teams [71]. We propose to further investigate this premise via a Markov decision process (MDP) to determine how teams should be dynamically assigned to work, isolate, or quarantine as

the number of infected physicians and the minimum capacity required for COVID-19 cases change simultaneously during an outbreak. In the following, the elements of the MDP are described.

### 3.3.1 Decision Epochs

At the end of each week, employees who have worked or are in isolation on that week are tested, and the ones whose test results are positive are quarantined. Since the tests are done every week, the decisions are made weekly.

### 3.3.2 The State Space

We assume that the number of teams is fixed and is denoted by  $I$ . As mentioned before, we considered teams of physicians that work in the same specialty. We define  $m_i^t$  as the total number of employees of team  $i$  at week  $t$ , and  $m^t = \sum_{i=1}^I m_i^t$ . We let the state of the system have the following structure:

$$\vec{S}^t = ((x_1^t, w_1^t, n_1^t, n1_1^t, n2_1^t, q_1^{n(t-1)}, P_1^{H(t-1)}), \dots, (x_I^t, w_I^t, n_I^t, n1_I^t, n2_I^t, q_I^{n(t-1)}, P_I^{H(t-1)}), P^t) \quad (3.1)$$

where  $x_i^t = \{0(Isolated), 1(Working)\}$  shows the status of team  $i$  at week  $t$ ,  $w_i^t$  is the number of weeks team  $i$  has been in the current status, and  $n_i^t$  is the number of employees of group  $i$  who have been in status  $x_i^t$  for  $w_i^t$  weeks. Moreover,  $nj_i$  is the number of employees in team  $i$  who are quarantined for  $j$  weeks. In each week, we also need to keep track of the number of employees who got a false negative result in previous week. The total number of employees of each group at each week is represented by  $m_i^t$ , where

$$m_i^t = n_i^t + n1_i^t + n2_i^t, \quad \forall i = 1, \dots, I \quad (3.2)$$

Further, we let  $P_i^{Ht}$  be the probability of positive test result of team  $i$  at week  $t$  and  $P^t$  be the this probability for the general population at week  $t$ . In this study, we consider a finite state space, thus we assume that the initial total number of employees in each team is finite. Furthermore, we assume that there is an upper bound on the number of weeks each team is allowed to be isolated or working, and Let  $w^1 < \infty$  and  $w^2 < \infty$  be the maximum number of weeks each team can work and be isolated, respectively. Since  $P_i^{Ht}$  and  $P^t$  are probabilities, we discretize them and assume that

they take values in finite sets such as  $P^t \in \{p_1, p_2, \dots, p_u\}$  and  $P_i^{Ht} \in \{p_1^{Ht}, p_2^{Ht}, \dots, p_{u_H}^{Ht}\}$ .

### 3.3.3 The Action Set

The actions that have to be taken at the end of each week determine the status of teams for the next week. Thus, we let the action set has the following structure:

$$A(\vec{s}^t) = \{a_i^t, i = 1, \dots, I\} \quad (3.3)$$

where  $a_i^t = \{0(Isolated), 1(Working)\}$  shows the action taken for team  $i$ . To guarantee a certain level of care capacity in each week, we assume that at minimum  $r(p^t)$  physicians are required to work.

$$\sum_{i=1}^I x_i^t n_i^t \geq r(p^t) \quad (3.4)$$

Note that the minimal capacity required depends on the infection probability in the general population to reflect the workload created by COVID-19. The following equalities define the set of allowable actions when a team has been working or isolated for the maximum number of weeks permitted.

$$A((x_i^t = 1, w_i^t, n^t, n1_i^t, n2_i^t, q_i^{n(t-1)}, p_i^{H(t-1)}), \dots) = (a_i^t = 0, \dots) \text{ if } w_i^t = w^1 \quad (3.5)$$

$$A((x_i^t = 0, w_i^t, n^t, n1_i^t, n2_i^t, q_i^{n(t-1)}, p_i^{H(t-1)}), \dots) = (a_i^t = 1, \dots) \text{ if } w_i^t = w^2 \quad (3.6)$$

### 3.3.4 Transitions

As mentioned above,  $P^t$  determines the positive test probability of the general population at week  $t$ . However, probability of getting infected is higher for healthcare employees because of being in direct contact with patients and their colleagues [89]. Let  $P_i^{Ht}$  be an estimation of the probability of getting a positive result in test for healthcare employees of team  $i$ . The value of  $P_i^{Ht}$  depends on the number of infected people who come to the healthcare facility, employees contact with them, and interaction of employees with each other. We define  $\rho_{ii'}$  as the probability that employees of team  $i$  contact employees of group  $i'$ , and  $\xi_i$  as the COVID-19 transmission probability of an employee of team  $i$  in case of contacting with another infected person. Moreover, we let  $\nu_i$  be

the probability that an employee of team  $i$  contacts patients. Thus,  $P_i^{Ht}$  can be defined as

$$P_i^{Ht} = p_j^{Ht}, \text{ if } p_{j-1}^{Ht} < P^t + x_i^t \left( \sum_{i'=1}^I x_{i'}^t \rho_{ii'} \xi_i \frac{y_{i'}^{1t}}{m_{i'}^t} + \nu_i \xi_i P^t \right) \leq p_j^{Ht} \text{ for } j \in \{1, \dots, u_H\} \quad (3.7)$$

and

$$P_i^{Ht} = p_{u_H}^{Ht}, \text{ if } p_{u_H}^{Ht} < P^t + x_i^t \left( \sum_{i'=1}^I x_{i'}^t \rho_{ii'} \xi_i \frac{y_{i'}^{1t}}{m_{i'}^t} + \nu_i \xi_i P^t \right) \quad (3.8)$$

where  $p_0^{Ht} = 0$  and  $y_i^{1t}$  is the number of working employees of team  $i$  who get a positive test result in week  $t$  which follows the binomial distribution:

$$y_i^{1t} \sim B(x_i^t(n_i^t - q_i^{n(t-1)}), P_i^{H(t-1)}), \quad \forall i = 1, \dots, I \quad (3.9)$$

In the above formulation,  $q_i^{n(t-1)}$  shows the number of infected employees of group  $i$  who got false negative result test at week  $t-1$ , where  $q_i^{n(t)} \sim B(y_i^{1t} + y_i^{2t}, p^n)$ , and  $p^n$  is the probability of getting false-negative result in the test. It is assumed that individuals who got a false negative test result in the previous week, show symptoms and/or are detected by the end of the current week. This assumption is consistent with empirical evidence that show symptomatic COVID-19 patients develop symptoms within 5 to 12 days [55]. Similarly, we let  $y_i^{2t}$  be the number of isolated employees of team  $i$  who get positive test result at week  $t$ .

$$y_i^{2t} \sim B((1 - x_i^t)(n_i^t - q_i^{n(t-1)}), P^{(t-1)}), \quad \forall i = 1, \dots, I \quad (3.10)$$

Let  $q_i^{1t}$  and  $q_i^{2t}$  be the number of employees of team  $i$  who passed away because of infection at week  $t$  after one and two weeks of quarantine, respectively. They follow Binomial distributions as follows:

$$q_i^{1t} \sim B(n1_i^t, \alpha^1), \quad \forall i = 1, \dots, I \quad (3.11)$$

$$q_i^{2t} \sim B(n2_i^t, \alpha^2), \quad \forall i = 1, \dots, I \quad (3.12)$$

where  $\alpha^1$  and  $\alpha^2$  are the mortality probability of employees in the first and second week of the disease. We also let  $r_i^t$  be the number of new arrivals at week  $t$ .  $r_i^t$  follows Poisson distribution with



rate  $\mu_i$ .

Thus, the evolution of elements of the state space is described as follows:

$$n1_i^{t+1} = y_i^{1t} + y_i^{2t} + q_i^{n(t-1)}, \quad \forall i = 1, \dots, I \quad (3.13)$$

$$n2_i^{t+1} = n1_i^t - q_i^{1t}, \quad \forall i = 1, \dots, I \quad (3.14)$$

$$n_i^{t+1} = (n_i^t + n1_i^t + n2_i^t + r_i^t - q_i^{1t} - q_i^{2t}) - n1_i^{t+1} - n2_i^{t+1}, \quad \forall i = 1, \dots, I \quad (3.15)$$

$$x_i^{t+1} = a_i^t, \quad \forall i = 1, \dots, I \quad (3.16)$$

$$w_i^{t+1} = \begin{cases} w_i^t + 1, & \text{if } x_i^t = a_i^t. \\ 1, & \text{otherwise.} \end{cases}, \quad \forall i = 1, \dots, I \quad (3.17)$$

The evolution of  $P^t$  is based on the real data released for the number of infected new cases diagnosed. Since our healthcare system of interest is located in South Carolina, we use the data for this area which is presented in Section 3.5.

### 3.3.5 The rewards

The objective is to schedule, isolate, or quarantine teams of healthcare employees to maximize the total expected discounted number of employees that are working to maintain capacity levels that can meet the patient demand for COVID-19 care. Thus, the immediate reward function has the following form:

$$R(\vec{s}^t, \vec{a}^t) = N^t \quad (3.18)$$

where

$$N^t = \sum_{i=1}^I \left( (n_i^t + n2_i^t + r_i^t - q_i^{1t} - q_i^{2t}) - y_i^{1t} - y_i^{2t} - q_i^{n(t-1)} \right)_{\{a_i^t=1\}} \quad (3.19)$$

### 3.3.6 Optimality equation

The value function  $v(\vec{s})$  corresponds to the total expected discounted reward (number of working employees) over the infinite horizon.

$$v(\vec{s}) = \max_{\vec{a} \in A_{\vec{s}}} \left\{ R(\vec{s}, \vec{a}) + \lambda \mathbb{E}(v(\vec{s}')) \right\}, \quad \forall \vec{s} \in S, \quad (3.20)$$

where  $\lambda \in [0, 1)$  is a discount factor. To find the optimal policy we need to solve equation (3.20). Since the state space and action set are finite, there exists a stationary optimal policy. However, the structure of the problem and size of the state space make a direct solution to (3.20) impractical. Therefore, we choose to develop an approximation approach that reduces the problem size.

## 3.4 Solution Approach

As mentioned in the previous section, due to the large size of the state space, it is impractical to find an exact solution to the MDP model. Thus, we employ model reduction techniques for computing an approximately optimal solution to the MDP. These techniques rely on finding a homogeneous partition of the state space where states in the same block of the partition transit with the same probability to each of the other blocks [26]. This partition makes a smaller MDP whose states are the blocks of the partition. We use the concept of  $\epsilon$ -homogeneous partition that is presented in [26]. The  $\epsilon$ -homogeneous partition allows states within the same block to transit with different probabilities to other blocks as long as the different probabilities are less than some  $\epsilon > 0$ . The  $\epsilon$ -homogeneous partitions are usually smaller than the homogeneous partitions [26]. In the following, we explain the  $\epsilon$ -homogeneous partition algorithm that is derived from [26].

**Definition 1** A partition  $P = \{B_1, \dots, B_n\}$  of the state space of an MDP  $M$  has the property of  $\epsilon$ -approximate stochastic bisimulation homogeneity with respect to  $M$  for  $\epsilon \in [0, 1]$  if and only if for each  $B_i, B_j \in P$ , for each  $\alpha \in A$ , and for each  $p, q \in B_i$ ,

$$|R(p) - R(q)| \in \epsilon, \quad (3.21)$$

$$\left| \sum_{r \in B_j} F_{pr}(\alpha) - \sum_{r \in B_j} F_{qr}(\alpha) \right| \leq \epsilon, \quad (3.22)$$

where  $A$  is set of actions,  $F$  assigns a probability to each state transition for each action, and  $R$  is

a reward function that maps each state to a real value. To construct an  $\epsilon$ -homogeneous partition, we first describe the relationship between every  $\epsilon$ -homogeneous partition and a particular simple partition based on immediate rewards.

**Definition 2** A partition  $P'$  is a refinement of a partition  $P$  if and only if each block of  $P'$  is a subset of some block of  $P$ ; in this case, we say that  $P$  is coarser than  $P'$ , and is a clustering of  $P'$ .

**Definition 3** The immediate reward partition is a partition in which two states,  $p$  and  $q$ , are in the same block if and only if they have the same reward.

**Definition 4** A partition  $P$  is  $\epsilon$ -uniform with respect to a function  $f : Q \rightarrow t$  if for every two states  $p$  and  $q$  in the same block of  $P$ ,  $|f(p) - f(q)| < \epsilon$ .

Every  $\epsilon$ -homogeneous partition is a refinement of some  $\epsilon$ -uniform clustering (with respect to reward) of the immediate reward partition [26]. The algorithm starts by constructing an  $\epsilon$ -uniform reward clustering  $P_0$  of the immediate reward partition. Then, this partition is refined by splitting blocks repeatedly to achieve  $\epsilon$ -homogeneity. To choose blocks for splitting, the following local property of the blocks of an  $\epsilon$ -homogeneous partition is used [26]:

**Definition 5** We say that a block  $C$  of a partition  $P$  is  $\epsilon$ -stable with respect to a block  $B$  if and only if for all actions  $\alpha$  and all states  $p \in C$  and  $q \in C$  we have  $|\sum_{r \in B} F_{pr}(\alpha) - \sum_{r \in B} F_{qr}(\alpha)| \leq \epsilon$ . We say that  $C$  is  $\epsilon$ -stable if  $C$  is  $\epsilon$ -stable with respect to every block of  $P$  and action in  $A$ .

According to the above definitions, a partition is  $\epsilon$ -homogeneous if and only if every block in the partition is  $\epsilon$ -stable. The algorithm checks each block for  $\epsilon$ -stability, splitting unstable blocks until there are no unstable blocks left to split [26]. When a block is unstable with respect to another block, we replace that by a set of sub-blocks that are  $\epsilon$ -stable with respect to the other blocks. Each  $\epsilon$ -homogeneous partition  $P$  of an MDP  $M$  induces a corresponding bounded parameter MDP (BMDP)  $M_P$ . The reward and transition ranges for blocks  $B$  and  $C$  of BMDP  $M_P$  and action  $\alpha$  are defined as:

$$\hat{R}(B) = [\min_{p \in B} R(p), \max_{p \in B} R(p)] \quad (3.23)$$

$$\hat{F}_{B,C}(\alpha) = [\min_{p \in B} \sum_{q \in C} F_{p,q}(\alpha), \max_{p \in B} \sum_{q \in C} F_{p,q}(\alpha)] \quad (3.24)$$

For a specific block, the reward value and transition probabilities are estimated by getting an average over all the states of that block.

Choosing a reasonable value for  $\epsilon$  is important in this algorithm. [68] proposed a way to equilibrate the value of  $\epsilon$  in equations 3.21 and 3.22. According to [68] formulation, two constant parameters  $C_p$  and  $C_r$  are defined and equations 3.21 and 3.22 are rewritten as

$$|R(p) - R(q)| \leq \epsilon C_r, \quad (3.25)$$

$$|\sum_{r \in B_j} F_{pr}(\alpha) - \sum_{r \in B_j} F_{qr}(\alpha)| \leq \epsilon C_p, \quad (3.26)$$

After using state reduction technique, the common methods for solving MDPs (policy iteration value iteration and linear programming) may be applied to solve the created BMDP.

## 3.5 Numerical Study and Numerical Results

In this section, we use our formulation and solution approach to obtain the optimal policy for staff scheduling under different scenarios. The selection of model parameters were informed by our discussions with Prisma Health. We also describe how the publicly available data on active COVID-19 cases in South Carolina can be used together with our model to analyze the effects of scheduling decisions under different pandemic conditions.

### 3.5.1 Input data

As mentioned in Section 3.3, estimation of  $P^t$  can be done based on the real data that is published for the number of new cases. Figure 3.1 shows the daily number of new detected cases since the first detected case's date up to July 1st, 2020 in South Carolina [3]. As it can be seen in Figure 3.1, the number of new cases has an increasing trend in the beginning, then for almost 2 months it keeps a stable trend. At the beginning of June 2020, the number of new cases again starts increasing rapidly. Since we are currently approaching a peak in South Carolina and there is no data available that shows the decreasing trend of data after passing the peak, we use the data from New York state to estimate the trend in the number of new infections in the subsequent months. In Figure 3.2 we show our forecast (after the red line) for the estimated number of new cases in South Carolina by investigating the similar trend in New York. This can provide us an estimation for the number positive cases toward the peak and after that until reaching a stable rate.

The value of  $P^t$  at each time period depends on its value in the previous period ( $P^{t-1}$ ) and

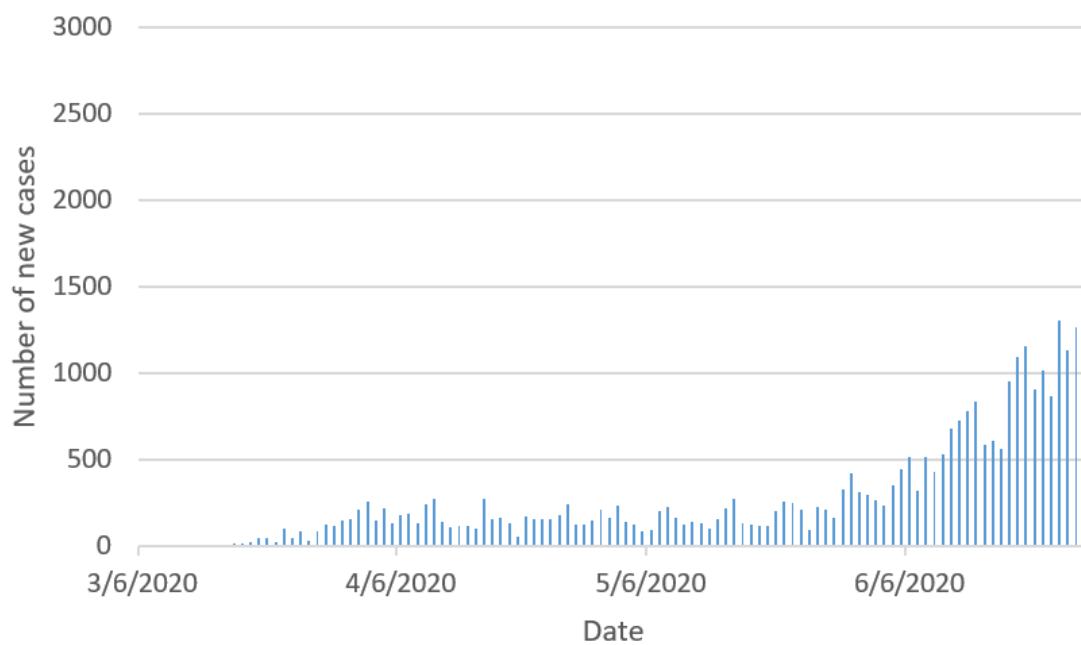


Figure 3.1: Daily new cases in South Carolina

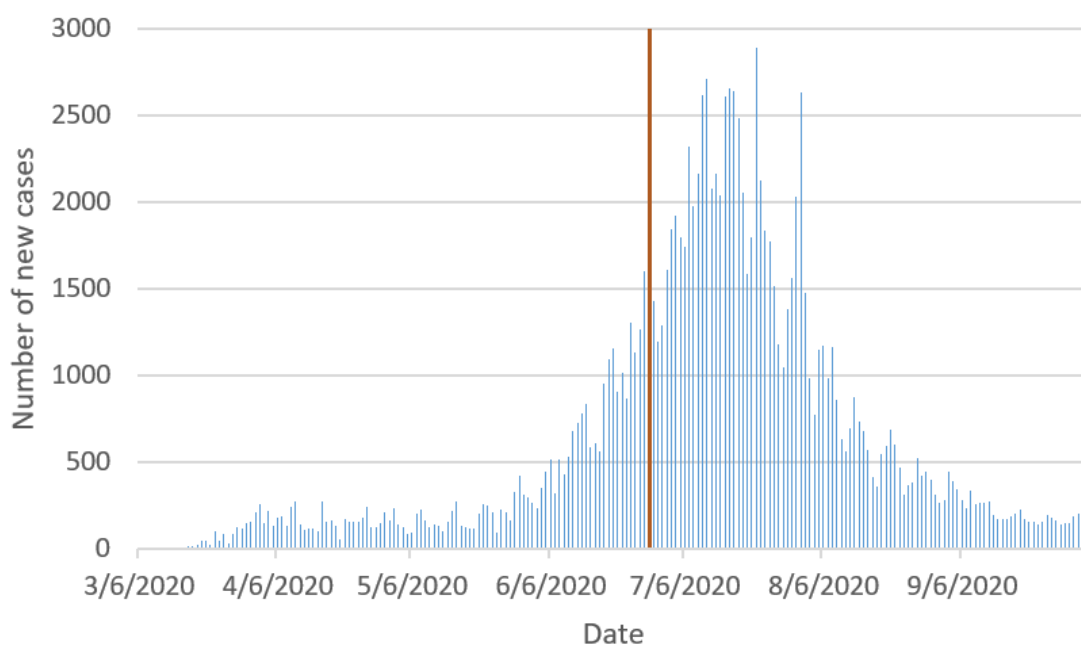


Figure 3.2: Forecast for the daily new cases in South Carolina using New York data

the public health measures taken in that period. Without a model that explains the relationship between the number of cases and the public health policies implemented, we are unable to estimate the value of  $P^t$  in such a way that truly shows the fluctuations over time. However, we can examine the behavior of our model under different infection rates to understand how the scheduling policies for healthcare employees should change through the different phases of the pandemic. To estimate the value of  $P^t$  at different points in time, we smoothed the data by applying 14-day moving average. The resulting data are shown in Figure 3.3.

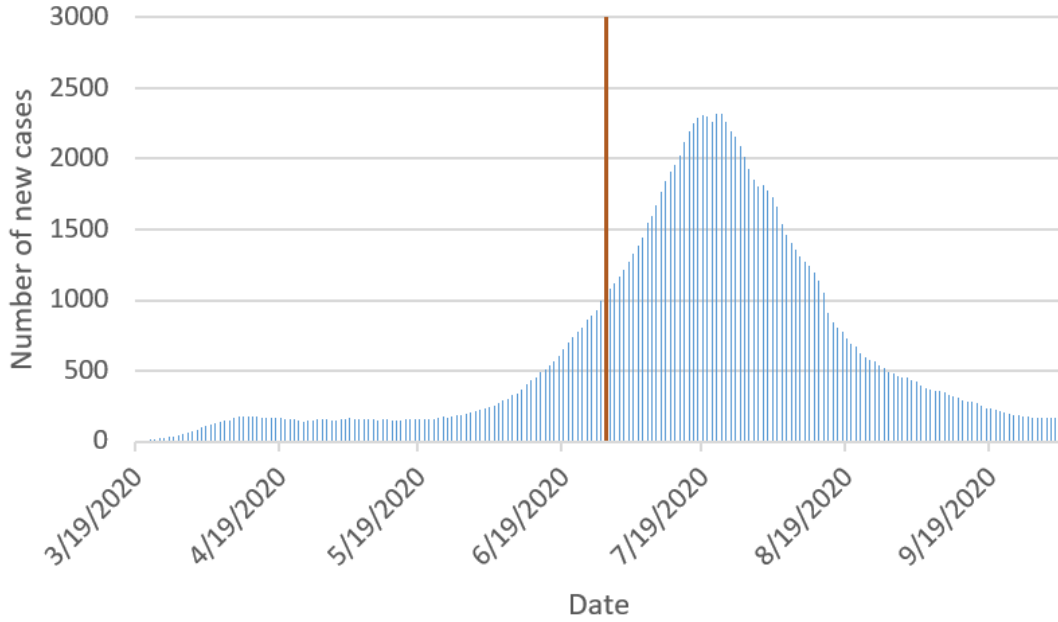


Figure 3.3: Smoothed data for the number of new cases

### 3.5.2 State reduction

We use state reduction technique presented in Section 1.4 to overcome the computational complexity presented by the model. We implemented the state reduction algorithm in C++. In the following, we define some scenarios and apply this algorithm with the corresponding set of parameters to observe how using this technique reduces size of the state space. To determine the value of  $\epsilon$  in this algorithm we set  $C_p = 1$  and  $C_r$  as the maximum possible value of the reward function which is 8 (the total number of employees). Thus,  $\epsilon$  can get any value within  $[1, 8]$ . We choose  $\epsilon = 6$  for these sets of experiments. The detailed description of the scenarios considered are

provided in Table 3.1. In Table 3.2, the size of each scenario and the size of reduced model after applying state reduction algorithm are shown.

Table 3.1: Estimation of parameters for each scenario

| Scenario | Parameters  |
|----------|---|
| 1        | $I = 2, M_1 = 4, M_2 = 4, w^1 = 3, w^2 = 3, p^t = 0.005, p^{Ht} \in \{0.125, 0.25, 0.375, 0.5\}$          |
| 2        | $I = 2, M_1 = 4, M_2 = 4, w^1 = 3, w^2 = 3, p^t = 0.01, p^{Ht} \in \{0.125, 0.25, 0.375, 0.5\}$           |
| 3        | $I = 2, M_1 = 4, M_2 = 4, w^1 = 3, w^2 = 3, p^t = 0.015, p^{Ht} \in \{0.125, 0.25, 0.375, 0.5\}$          |
| 4        | $I = 2, M_1 = 4, M_2 = 4, w^1 = 3, w^2 = 3, p^t = 0.005, p^{Ht} \in \{0.1, 0.2, 0.3, 0.4, 0.5\}$          |
| 5        | $I = 2, M_1 = 4, M_2 = 4, w^1 = 4, w^2 = 4, p^t = 0.005, p^{Ht} \in \{0.125, 0.25, 0.375, 0.5\}$          |
| 6        | $I = 3, M_1 = 3, M_2 = 3, M_3 = 2, p^t = 0.005, w^1 = 3, w^2 = 3, p^{Ht} \in \{0.125, 0.25, 0.375, 0.5\}$ |

Table 3.2: Comparison of the original and reduced problems' sizes

| Scenario | State space size in original problem | State space size in reduced problem |
|----------|--------------------------------------|-------------------------------------|
| 1        | 230,400                              | 4750                                |
| 2        | 230,400                              | 5112                                |
| 3        | 230,400                              | 5298                                |
| 4        | 360,000                              | 8013                                |
| 5        | 518,400                              | 18,780                              |
| 6        | 65,523,600                           | 3,950,324                           |

As it is shown in Table 3.2, the state reduction algorithm has a huge impact on decreasing the size of the state space. In particular, the number states is reduced by 95% to 98% in the scenarios considered. This sizable reduction means that the reduced problem can be solved by using standard methods in the literature. In the next section, we work with the reduced state space and compute an approximate optimal policy using the policy iteration algorithm.

### 3.5.3 Policy iteration to solve the reduced model

After obtaining the reduce model, we use policy iteration algorithm which is a common method of solving MDPs to solve the resulting problem [72]. The policy iteration algorithm is implemented in C++. To observe how the difference in infection probability for the general population and the number of employee teams affect the optimal policy, we solve scenarios 1, 2, 3 and 6. Op-

timal policy for each scenario can be analyzed from different view points since the state space has multiple elements and changes in the optimal actions along different dimensions of the state space can be observed. In this section, we present some of these observations that we think may be more insightful.

The value of  $p^t$  in scenarios 1 and 6 assumed to be constant and equal to 0.005. This value is estimated based on the smoothed number of new cases from May 2<sup>th</sup>, 2020 to June 3<sup>rd</sup>, 2020 and the total number of susceptible individuals living in South Carolina. Note that in Figure 3.3, the infection rate during this time period is relatively the stable. In scenarios 2 and 3, we investigate the impact of larger  $p^t$  in the optimal policy and consider  $p^t = 0.01$  and  $p^t = 0.015$  in scenarios 2 and 3, respectively. In Figure 3.4, we show these three values of  $p^t$  and their corresponding dates.

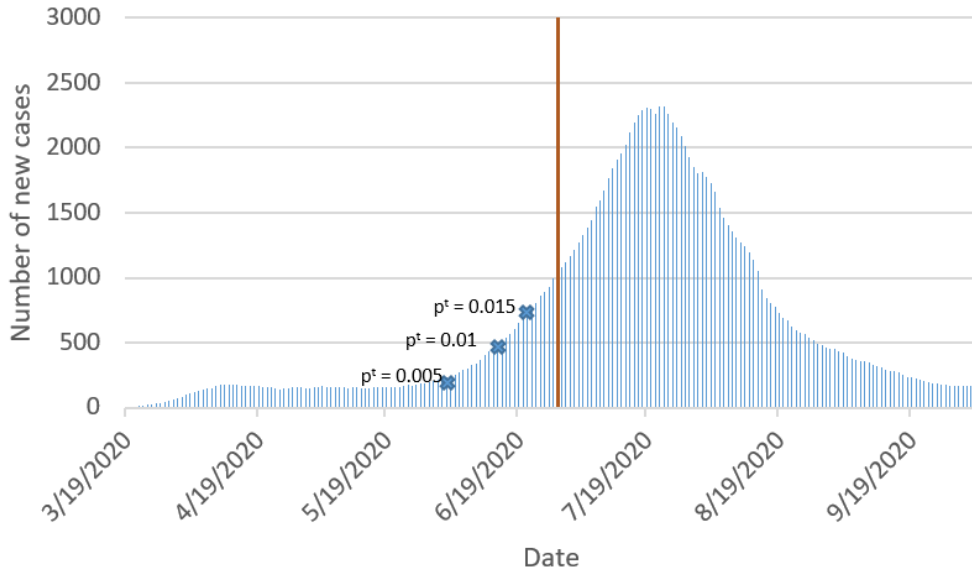


Figure 3.4: Smoothed data for the number of new cases with specified  $p^t$  values

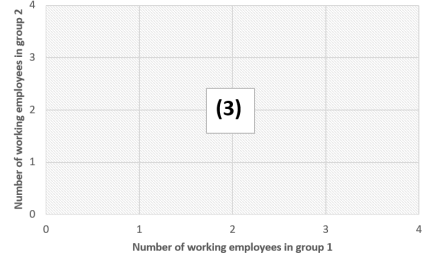
The figures 3.5, 3.6, 3.7 and 3.8 show the optimal action with respect to the number of working employees in each group. In each graph, one or multiple actions that are denoted by (1), (2) and (3) are depicted. Action (1) corresponds to only team 1 working, action (2) corresponds to only team 2 working, and Action (3) corresponds to both team 1 and 2 working. Note that in scenario 6 the number of teams is increased to 3. To make the presentation of our results simpler in Figure 3.8, optimal actions are depicted only for the set of states where the number of working employees in group 3 is fixed at 2. Furthermore, we only depict the optimal actions in states where



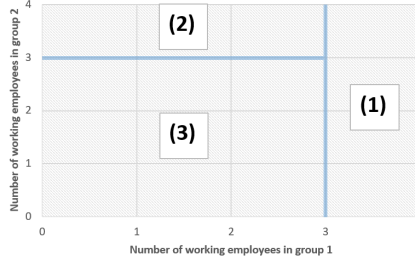
the values of  $w_1^t$  and  $w_2^t$  is one, which means both groups have been already working for one week.



(a)  $P^{Ht} = 0.125$



(b)  $P^{Ht} = 0.25$

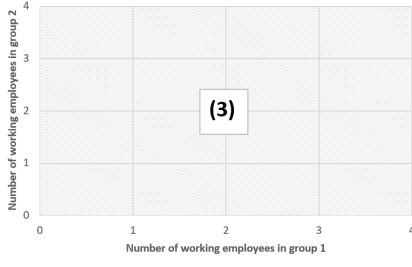


(c)  $P^{Ht} = 0.375$

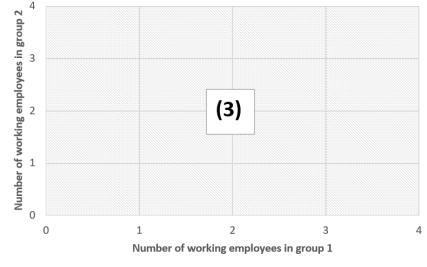
Figure 3.5: Optimal policy for scenario 1 when  $w_1^t = 1$ ,  $w_2^t = 1$

Figure 3.5 shows the optimal policy for scenario 1. Figures 3.5a, 3.5b, and 3.5c are the cases in which the value of  $P^{Ht}$  is 0.125, 0.25 and 0.375, respectively. Our results show that, when  $P^{Ht}$  is 0.125 or 0.25, the optimal action has both teams working. However, when  $P^{Ht}$  increases to 0.375, depending on the number of employees who are not in quarantine in each team, it can be optimal to isolate one of the teams. Similarly, Figure 3.6 shows the optimal policy for scenario 2. As it can be seen in these graphs, the optimal policy does not change when  $p^t$  is increased from 0.005 in scenario 1 to 0.01 in scenario 2. However, when the infection probability for the general population is higher at  $p^t = 0.015$  in scenario 3, it is not optimal for both teams to work and the optimal action will depend on the number of employees who are not quarantined in each team.

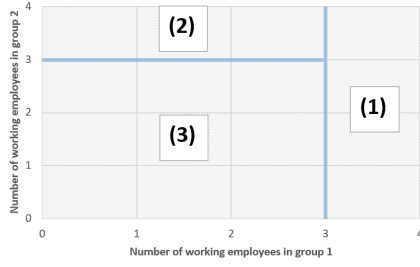
The optimal policy for scenario 6 where the number of groups is increased to 3 is shown in Figure 3.8. Our results show that for teams 1 and 2, when  $P^{Ht} = 0.125$  the optimal action will be to work. However, when the value of  $P^{Ht}$  is increased to 0.25 and 0.375, the optimal action is again dependent on the number of healthy employees in each group.



(a)  $P^{Ht} = 0.125$

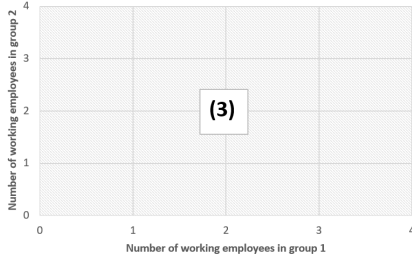


(b)  $P^{Ht} = 0.25$

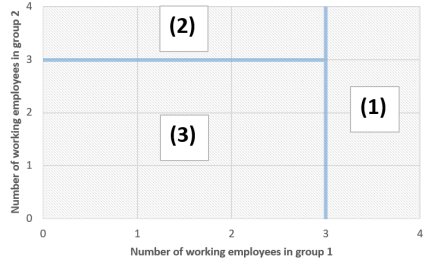


(c)  $P^{Ht} = 0.375$

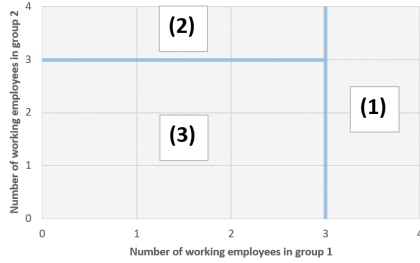
Figure 3.6: Optimal policy for scenario 2 when  $w_1^t = 1$ ,  $w_2^t = 1$



(a)  $P^{Ht} = 0.125$

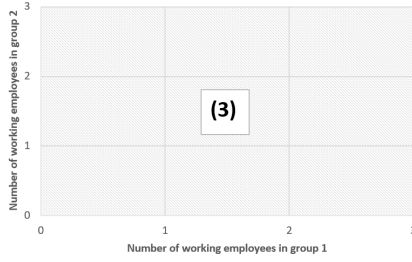


(b)  $P^{Ht} = 0.25$

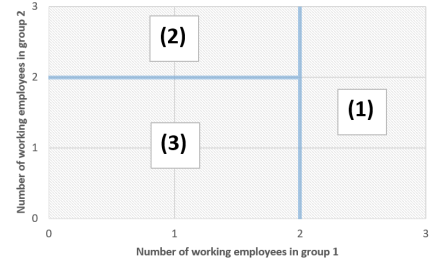


(c)  $P^{Ht} = 0.375$

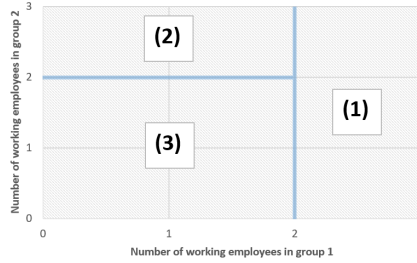
Figure 3.7: Optimal policy for scenario 3 when  $w_1^t = 1$ ,  $w_2^t = 1$



(a)  $P^{Ht} = 0.125$



(b)  $P^{Ht} = 0.25$



(c)  $P^{Ht} = 0.375$

Figure 3.8: Optimal policy for scenario 6 when  $w_1^t = 1$ ,  $w_2^t = 1$ ,  $n_3^t = 2$

To evaluate the performance of our solution methodology, we compare the total discounted reward achieved by the approximate optimal policy to that of a benchmark policy that has all teams working without isolation. Note that under this policy infected employees are still quarantined for two weeks. These results are shown in Table 3.3 for all the scenarios considered.

Table 3.3: Comparison of the approximate optimal policy versus the benchmark policy

| Scenario | Optimal policy | Benchmark - no isolation |
|----------|----------------|--------------------------|
| 1        | 9.6            | 8.9                      |
| 2        | 9.5            | 8.6                      |
| 3        | 9.2            | 8.5                      |
| 4        | 10.3           | 9.6                      |
| 5        | 10.9           | 9.2                      |
| 6        | 10.5           | 10                       |

Our results show that in all scenarios, the approximate optimal policy leads to an improvement over scheduling all teams to work at all times. The improvement achieved ranges from 5% to 18% for different scenarios, which can be of critical importance when the healthcare system is under

strain due to increased demand. In particular, the results show that when the infection probability increases moderately (in scenarios 2 and 3), the number of employees that are assigned to work under the approximate optimal policy decreases slightly, but the total discounted number of working employees decrease faster under the benchmark policy due to increased number of employees who have to quarantine for two weeks. Moreover, when the allowable number of weeks that each group can work in a row increases (scenario 5), the model assigns more employees to work and that increases the value function for this scenario compare to scenario 1. Finally, we observe that when the number of teams formed is increased (for the same total number of employees), the relative performance of the benchmark policy is improved due to the decreased interaction among the healthcare employees.

### 3.6 Conclusion

In this study, we introduce a dynamic model for scheduling of healthcare employees in the COVID-19 pandemic. Our model can be applied to study other infectious disease outbreaks as well. Previous studies show that organizing healthcare employees into teams during the COVID-19 pandemic and scheduling them based on as teams decreases the probability of contact with infected individuals and hence decreases the probability of getting infected. We employ this idea and build a MDP model in which an optimal policy that determines how these teams should be scheduled to work in long run to maximize the utilization of the personnel is sought. Since the size of the state space is large, before using policy iteration to solve the MDP model, we apply a state space reduction technique to decrease the size of the state space. That makes the problem easier to solve. By obtaining the approximate optimal policy under different scenarios and analyzing its performance, we show that voluntary intermittent isolation of some teams (as opposed to having all teams work continuously) can be beneficial to maximize the number of workable-physician-days while the limiting the risk of infection for physicians. This strategy is particularly useful when the COVID-19 infection rate in the general population is increasing.

This study can be extended in multiple ways. In our numerical experiments, the value of infection probability in this study is assumed to be constant under each scenario considered. However, the real data shows a more complicated trend in the value of  $p^t$ . Thus, we can extend this study by considering a function for infection probability. Additionally, using a finite horizon MDP, it might be a possible to update  $p^t$  over time to account for increasing or decreasing trends.

There are other clustering algorithms to reduce the state space of MDPs in the literature, and furthermore, even in the method we used, a strategic selection of states to cluster around could possibly be beneficial. So, one of the possible extensions of this study is considering these algorithms and compare the results with the method we used in this study.

## Chapter 4

# Contributions and Future Research

In this dissertation, we study three problems related to healthcare operations management by which we demonstrate how using mathematical optimization in novel ways can improve healthcare systems. Importance of mathematical optimization in decision making in healthcare systems has been subject of many previous studies. Our contribution to this existing literature lies in the study of newly developed concepts and models in areas of capacity allocation, scheduling, and infection detection and prevention among healthcare workers. In these projects, we collaborated with Prisma Health in South Carolina. Thereupon, this thesis consists of three tactical and operational problems in healthcare systems, and mathematical optimization techniques are used to address and solve them.

In the first study, we propose a novel model for scheduling of patients in a diagnostic clinic, and we show how our model improves waiting time of patients to receive appointments, and how the utilization of the system is improved. This model is called “postponement model” since our model intentionally postpones the scheduling of some lower priority patients to reserve appointments for probable higher priority patients. We believe that our study is the first in the literature that applies this concept in a scheduling model. This study can be extended in multiple directions. For example, postponing the acceptances may increase the possibility of no-shows. Thus, effects of postponement on no-show behavior and patient preferences can be part of future work. Furthermore, we assume that the duration of visits are constant and identical for each type of patient. Thus, the number of patients that can be seen each day is a fixed number in our model. To generalize these assumptions, uncertain or patient type-dependent service times can be considered.

In the second study, we propose a dynamic TB screening Model for healthcare employees which minimizes the total cost due to screening and infections in the healthcare facility. As far as we are aware, our study is the first in the literature that uses mathematical modeling in organizing TB screening for healthcare employees. TB screening is usually annually in healthcare facilities. Our analysis represent that yearly testing of all employees regardless of infection risk is not cost effective, and similarly low infection rates can be achieved via strategic intermittent testing. One of the possible extensions to our model is to differentiate between active and latent TB and consider its impact on the optimal policy. We believe that there are some parameters in our formulation that their value affect the results such as percentage of TB infected patients who visit the healthcare facility, probability of getting infected by contacting a sick person, and the probability of contacting with infected individuals. Performing sensitivity analysis tasks to observe how changing these parameters affect the optimal policy can also be part of future work in this area.

Finally, the focus of the third problem is managing the healthcare workforce during the COVID-19 pandemic to maintain the high utilization of the system while controlling the infection rate among the employees. Our study is the first in the literature that optimizes scheduling of teams of healthcare workers during a pandemic. The results of solving our model give potential policies to be applied to increase the number of working employees during COVID-19 pandemic while we consider the possibility of getting infected by working. In the experiments we performed in this study, we assume that the COVID-19 infection probability is constant for the general population. However, the infection trends evidenced by data over time shows that it may be constant, increasing or decreasing depending on which stage of the pandemic the patient population is. An extension to this study is considering a function for infection probability. By considering a finite horizon model, one may be able to update this probability over time. There are other clustering algorithms to reduce the state space of MDPs in the literature, and furthermore, even in the method we used, a strategic selection of states to cluster around could possibly be beneficial. So, one of the possible extensions of this study is considering these algorithms and compare the results with the method we used in this study.

# Appendices



## Appendix A

Table 1: Parameters values for the postponement model

| Parameter   | Value | Parameter   | Value | Parameter   | Value |
|-------------|-------|-------------|-------|-------------|-------|
| $\lambda^E$ | 50    | $b_{1,2,4}$ | 0.75  | $b_{2,2,5}$ | 2     |
| $\lambda^I$ | 20    | $b_{1,2,5}$ | 0.75  | $b_{2,2,6}$ | 2.5   |
| $\lambda_1$ | 80    | $b_{1,2,6}$ | 1.25  | $b_{2,2,7}$ | 2.5   |
| $\lambda_2$ | 40    | $b_{1,2,7}$ | 1.25  | $b_{2,3,1}$ | 3     |
| $c^E$       | 45    | $b_{1,3,1}$ | 1.5   | $b_{2,3,2}$ | 4     |
| $c^I$       | 30    | $b_{1,3,2}$ | 2.5   | $b_{2,3,3}$ | 4     |
| $c_{1,1}^O$ | 8     | $b_{1,3,3}$ | 2.5   | $b_{2,3,4}$ | 5     |
| $c_{1,2}^O$ | 12    | $b_{1,3,4}$ | 3.5   | $b_{2,3,5}$ | 5     |
| $c_{1,3}^O$ | 16    | $b_{1,3,5}$ | 3.5   | $b_{2,3,6}$ | 6     |
| $c_{2,1}^O$ | 14    | $b_{1,3,6}$ | 4.5   | $b_{2,3,7}$ | 6     |
| $c_{2,2}^O$ | 18    | $b_{1,3,7}$ | 4.5   | $a_{1,1}$   | 1     |
| $c_{2,3}^O$ | 22    | $b_{2,1,1}$ | 0     | $a_{1,2}$   | 1.25  |
| $b_{1,1,1}$ | 0     | $b_{2,1,2}$ | 0.5   | $a_{1,3}$   | 4.5   |
| $b_{1,1,2}$ | 0     | $b_{2,1,3}$ | 0.5   | $a_{2,1}$   | 1.5   |
| $b_{1,1,3}$ | 0.5   | $b_{2,1,4}$ | 1     | $a_{2,2}$   | 2.5   |
| $b_{1,1,4}$ | 0.5   | $b_{2,1,5}$ | 1     | $a_{2,3}$   | 6     |
| $b_{1,1,5}$ | 0.5   | $b_{2,1,6}$ | 1.5   | $w_{1,1}$   | 1     |
| $b_{1,1,6}$ | 1     | $b_{2,1,7}$ | 1.5   | $w_{1,2}$   | 0.95  |
| $b_{1,1,7}$ | 1     | $b_{2,2,1}$ | 1     | $w_{1,3}$   | 0.9   |
| $b_{1,2,1}$ | 0.25  | $b_{2,2,2}$ | 1.5   | $w_{2,1}$   | 0.95  |
| $b_{1,2,2}$ | 0.25  | $b_{2,2,3}$ | 1.5   | $w_{2,2}$   | 0.9   |
| $b_{1,2,3}$ | 0.75  | $b_{2,2,4}$ | 2     | $w_{2,3}$   | 0.85  |

Table 2: Estimated parameters in no-postponement model

| Parameter   | Estimation      | Parameter | Estimation | Parameter | Estimation |
|-------------|-----------------|-----------|------------|-----------|------------|
| $\lambda^E$ | 50              | $b_{1,1}$ | 0          | $b_{2,2}$ | 0.25       |
| $\lambda^I$ | 20              | $b_{1,2}$ | 0          | $b_{2,3}$ | 0.25       |
| $\lambda_1$ | $\frac{80}{18}$ | $b_{1,3}$ | 0.25       | $b_{2,4}$ | 0.5        |
| $\lambda_2$ | $\frac{40}{18}$ | $b_{1,4}$ | 0.25       | $b_{2,5}$ | 0.5        |
| $c^E$       | 45              | $b_{1,5}$ | 0.25       | $b_{2,6}$ | 0.75       |
| $c^I$       | 30              | $b_{1,6}$ | 0.5        | $b_{2,7}$ | 0.75       |
| $c_1^O$     | 6               | $b_{1,7}$ | 0.5        |           |            |
| $c_2^O$     | 10              | $b_{2,1}$ | 0          |           |            |

## Appendix B

Table 3: Parameters values

| Parameter      | Value | Parameter      | Value | Parameter  | Value    |
|----------------|-------|----------------|-------|------------|----------|
| $\lambda_{11}$ | 4     | $\nu_{23}$     | 1     | $p_{s1}^p$ | 0.6      |
| $\lambda_{12}$ | 14    | $\nu_{31}$     | 0.75  | $p_{s2}^p$ | 0.27     |
| $\lambda_{13}$ | 10    | $\nu_{32}$     | 0.5   | $p_{s3}^p$ | 0.27     |
| $\lambda_{21}$ | 15    | $\nu_{33}$     | 1     | $p_{s1}^n$ | 0.04     |
| $\lambda_{22}$ | 50    | $\rho_{11,11}$ | 1     | $p_{s2}^n$ | 0.04     |
| $\lambda_{23}$ | 35    | $\rho_{12,12}$ | 1     | $p_{s3}^n$ | 0.04     |
| $\lambda_{31}$ | 4     | $\rho_{13,13}$ | 1     | $p_{b1}^p$ | 0.176    |
| $\lambda_{32}$ | 14    | $\rho_{21,21}$ | 1     | $p_{b2}^p$ | 0.176    |
| $\lambda_{33}$ | 10    | $\rho_{22,22}$ | 1     | $p_{b3}^p$ | 0.176    |
| $p_{11}^l$     | 0.15  | $\rho_{23,23}$ | 1     | $p_{b1}^n$ | 0.008    |
| $p_{12}^l$     | 0.15  | $\rho_{31,31}$ | 1     | $p_{b2}^n$ | 0.008    |
| $p_{13}^l$     | 0.15  | $\rho_{32,32}$ | 1     | $p_{b3}^n$ | 0.008    |
| $p_{21}^l$     | 0.15  | $\rho_{33,33}$ | 1     | $c^b$      | 45 \$    |
| $p_{22}^l$     | 0.15  | $\xi_{11}$     | 0.05  | $c^s$      | 8 \$     |
| $p_{23}^l$     | 0.15  | $\xi_{12}$     | 0.22  | $c^x$      | 100 \$   |
| $p_{31}^l$     | 0.15  | $\xi_{13}$     | 0.22  | $c_1^l$    | 150 \$/h |
| $p_{32}^l$     | 0.15  | $\xi_{21}$     | 0.05  | $c_2^l$    | 30 \$/h  |
| $p_{33}^l$     | 0.15  | $\xi_{22}$     | 0.22  | $c_3^l$    | 29 \$/h  |
| $\nu_{11}$     | 1     | $\xi_{23}$     | 0.22  | $c_1^u$    | 5000 \$  |
| $\nu_{12}$     | 1     | $\xi_{31}$     | 0.05  | $c_2^u$    | 1000 \$  |
| $\nu_{13}$     | 1     | $\xi_{32}$     | 0.22  | $c_3^u$    | 1000 \$  |
| $\nu_{21}$     | 1     | $\xi_{33}$     | 0.22  |            |          |
| $\nu_{22}$     | 1     | $\beta$        | 0.1   |            |          |

# Bibliography

- [1] *Chest X-Ray Cost and Procedure Information*, 2020. <https://www.newchoicehealth.com/procedures/chest-x-ray>.
- [2] *TB Blood Test*, 2020. <https://www.mdsave.com/procedures/tb-blood-test/d582fdc4>.
- [3] *COVID-19 OUTBREAK*, 2020-06-27. <https://www.postandcourier.com/health/covid19/>.
- [4] Amir Ahmadi-Javid, Zahra Jalali, and Kenneth J Klassen. Outpatient appointment systems in healthcare: A review of optimization studies. *European Journal of Operational Research*, 258(1):3–34, 2017.
- [5] Niyaz Ahmed and Seyed E Hasnain. Molecular epidemiology of tuberculosis in india: Moving forward with a systems biology approach. *Tuberculosis*, 91(5):407–413, 2011.
- [6] Shabbir Ahmed, Alexander Shapiro, and Er Shapiro. The sample average approximation method for stochastic programs with integer recourse. *Submitted for publication*, pages 1–24, 2002.
- [7] Shabbir Ahmed, Mohit Tawarmalani, and Nikolaos V Sahinidis. A finite branch-and-bound algorithm for two-stage stochastic integer programs. *Mathematical Programming*, 100(2):355–377, 2004.
- [8] Raha Akhavan-Tabatabaei, Diana Marcela Sánchez, and Thomas G Yeung. A markov decision process model for cervical cancer screening policies in colombia. *Medical Decision Making*, 37(2):196–211, 2017.
- [9] Faezeh Akhavizadegan, Javad Ansarifard, and Fariborz Jolai. A novel approach to determine a tactical and operational decision for dynamic appointment scheduling at nuclear medical center. *Computers & Operations Research*, 78:267–277, 2017.
- [10] Oguzhan Alagoz, Cindy L Bryce, Steven Shechter, Andrew Schaefer, Chung-Chou H Chang, Derek C Angus, and Mark S Roberts. Incorporating biological natural history in simulation models: empirical estimates of the progression of end-stage liver disease. *Medical Decision Making*, 25(6):620–632, 2005.
- [11] Oguzhan Alagoz, Lisa M Maillart, Andrew J Schaefer, and Mark S Roberts. The optimal timing of living-donor liver transplantation. *Management Science*, 50(10):1420–1430, 2004.
- [12] Oguzhan Alagoz, Lisa M Maillart, Andrew J Schaefer, and Mark S Roberts. Choosing among living-donor and cadaveric livers. *Management Science*, 53(11):1702–1715, 2007.
- [13] Emanuele Amodio, Giovanna Anastasi, Maria Grazia Laura Marsala, Maria Valeria Torregrossa, Nino Romano, and Alberto Firenze. Vaccination against the 2009 pandemic influenza a (h1n1) among healthcare workers in the major teaching hospital of sicily (italy). *Vaccine*, 29(7):1408–1412, 2011.

- [14] Turgay Ayer, Oguzhan Alagoz, and Natasha K Stout. Or forum—a pomdp approach to personalize mammography screening decisions. *Operations Research*, 60(5):1019–1034, 2012.
- [15] Hari Balasubramanian, Ana Muriel, Asli Ozen, Liang Wang, Xiaoling Gao, and Jan Hippchen. Capacity allocation and flexibility in primary care. In *Handbook of healthcare operations management*, pages 205–228. Springer, 2013.
- [16] Achal Bassamboo, J Michael Harrison, and Assaf Zeevi. Dynamic routing and admission control in high-volume service systems: Asymptotic analysis via multi-scale fluid limits. *Queueing Systems*, 51(3-4):249–285, 2005.
- [17] Bjorn P Berg, Brian T Denton, S Ayca Erdogan, Thomas Rohleder, and Todd Huschka. Optimal booking and scheduling in outpatient procedure centers. *Computers & Operations Research*, 50:24–37, 2014.
- [18] Dimitris Bertsimas and Ioana Popescu. Revenue management in a dynamic network environment. *Transportation science*, 37(3):257–277, 2003.
- [19] Papiya Bhattacharjee and Pradip Kumar Ray. Simulation modelling and analysis of appointment system performance for multiple classes of patients in a hospital: a case study. *Operations Research for Health Care*, 8:71–84, 2016.
- [20] Mucahit Cevik, Turgay Ayer, Oguzhan Alagoz, and Brian L Sprague. Analysis of mammography screening policies under resource constraints. *Production and Operations Management*, 27(5):949–972, 2018.
- [21] Jagpreet Chhatwal, Oguzhan Alagoz, and Elizabeth S Burnside. Optimal breast biopsy decision-making based on mammographic features and demographic factors. *Operations research*, 58(6):1577–1591, 2010.
- [22] Josette SY Chor, Surinder K Pada, Iain Stephenson, William B Goggins, Paul A Tambyah, Tristan William Clarke, Mariejo Medina, Nelson Lee, Ting Fun Leung, Karry LK Ngai, et al. Seasonal influenza vaccination predicts pandemic h1n1 vaccination uptake among healthcare workers in three countries. *Vaccine*, 29(43):7364–7369, 2011.
- [23] Philip Cooley, Bruce Y Lee, Shawn Brown, James Cajka, Bernadette Chasteen, Laxminarayana Ganapathi, James H Stark, William D Wheaton, Diane K Wagener, and Donald S Burke. Protecting health care workers: a pandemic simulation based on allegheny county. *Influenza and other respiratory viruses*, 4(2):61–72, 2010.
- [24] Natasha S Crowcroft, CE Roth, Bernard J Cohen, and Elizabeth Miller. Guidance for control of parvovirus b19 infection in healthcare settings and the community. *Journal of Public Health*, 21(4):439–446, 1999.
- [25] Marie A De Perio, Joel Tsevat, Gary A Roselle, Stephen M Kralovic, and Mark H Eckman. Cost-effectiveness of interferon gamma release assays vs tuberculin skin tests in health care workers. *Archives of internal medicine*, 169(2):179–187, 2009.
- [26] Thomas L Dean, Robert Givan, and Sonia Leach. Model reduction techniques for computing approximately optimal solutions for markov decision processes. *arXiv preprint arXiv:1302.1533*, 2013.
- [27] Jivan Deglise-Hawkinson, Jonathan E Helm, Todd Huschka, David L Kaufman, and Mark P Van Oyen. A capacity allocation planning model for integrated care and access management. *Production and operations management*, 27(12):2270–2290, 2018.

- [28] Susan E Dorman, Robert Belknap, Edward A Graviss, Randall Reves, Neil Schluger, Paul Wein-furter, Yaping Wang, Wendy Cronin, Yael Hirsch-Moverman, Larry D Teeter, et al. Interferon- $\gamma$  release assays and tuberculin skin testing for diagnosis of latent tuberculosis infection in health-care workers in the united states. *American journal of respiratory and critical care medicine*, 189(1):77–87, 2014.
- [29] JS Duchin, JA Jereb, CM Nolan, P Smith, and IM Onorato. Comparison of sensitivities to two commercially available tuberculin skin test reagents in persons with recent tuberculosis. *Clinical infectious diseases*, 25(3):661–663, 1997.
- [30] Maryam Eghbali-Zarch, Reza Tavakkoli-Moghaddam, Fatemeh Esfahanian, Amir Azaron, and Mohammad Mehdi Sepehri. A markov decision process for modeling adverse drug reactions in medication treatment of type 2 diabetes. *Proceedings of the Institution of Mechanical Engineers, Part H: Journal of Engineering in Medicine*, page 0954411919853394, 2019.
- [31] S Ayca Erdogan and Brian Denton. Dynamic appointment scheduling of a stochastic server with uncertain demand. *INFORMS Journal on Computing*, 25(1):116–132, 2013.
- [32] S Ayca Erdogan, Alexander Gose, and Brian T Denton. Online appointment sequencing and scheduling. *IIE Transactions*, 47(11):1267–1286, 2015.
- [33] Jacob Feldman, Nan Liu, Huseyin Topaloglu, and Serhan Ziya. Appointment scheduling under patient preference and no-show behavior. *Operations Research*, 62(4):794–811, 2014.
- [34] Martina Ferioli, Cecilia Cisternino, Valentina Leo, Lara Pisani, Paolo Palange, and Stefano Nava. Protecting healthcare workers from sars-cov-2 infection: practical indications. *European Respiratory Review*, 29(155), 2020.
- [35] Sarah A Foster-Chang, Mary L Manning, and Laura Chandler. Tuberculosis screening of new hospital employees: Compliance, clearance to work time, and cost using tuberculin skin test and interferon-gamma release assays. *Workplace health & safety*, 62(11):460–467, 2014.
- [36] Long Gao, Susan H Xu, and Michael O Ball. Managing an available-to-promise assembly system with dynamic short-term pseudo-order forecast. *Management Science*, 58(4):770–790, 2012.
- [37] Haileyesus Getahun, Alberto Matteelli, Richard E Chaisson, and Mario Raviglione. Latent mycobacterium tuberculosis infection. *New England Journal of Medicine*, 372(22):2127–2135, 2015.
- [38] Linda V Green, Sergei Savin, and Ben Wang. Managing patient service in a diagnostic medical facility. *Operations Research*, 54(1):11–25, 2006.
- [39] Diwakar Gupta and Brian Denton. Appointment scheduling in health care: Challenges and opportunities. *IIE transactions*, 40(9):800–819, 2008.
- [40] Riley Hazard, Kyle B Enfield, Darla J Low, Eve T Giannetta, and Costi D Sifri. Hidden reservoir: An outbreak of tuberculosis in hospital employees with no patient contact. *infection control & hospital epidemiology*, 37(9):1111–1113, 2016.
- [41] Elizabeth Hunter, Brian Mac Namee, and John D Kelleher. A taxonomy for agent-based models in human infectious disease epidemiology. *Journal of Artificial Societies and Social Simulation*, 20(3), 2017.
- [42] Paul A Jensen, Lauren A Lambert, Michael F Iademarco, and Renee Ridzon. Guidelines for preventing the transmission of mycobacterium tuberculosis in health-care settings, 2005. 2005.

- [43] Ruiwei Jiang, Siqian Shen, and Yiling Zhang. Integer programming approaches for appointment scheduling with random no-shows and service durations. *Operations Research*, 65(6):1638–1656, 2017.
- [44] Günter Kampf, Daniel Todt, Stephanie Pfaender, and Eike Steinmann. Persistence of coronaviruses on inanimate surfaces and their inactivation with biocidal agents. *Journal of Hospital Infection*, 104(3):246–251, 2020.
- [45] Keumseok Kang, J George Shanthikumar, and Kemal Altinkemer. Postponable acceptance and assignment: A stochastic dynamic programming approach. *Manufacturing & Service Operations Management*, 18(4):493–508, 2016.
- [46] Edward H Kaplan. Probability models of needle exchange. *Operations Research*, 43(4):558–569, 1995.
- [47] Ben Kesling and Dion Nissenbaum. VA Goal to Slash Wait Times Was ‘Unrealistic’, Aide Said. *The Wall Street Journal (May 23)*. Available at <https://www.wsj.com/articles/SB10001424052702303749904579580473122138420> (accessed date August 27, 2017), 2014.
- [48] Kanwal Fatima Khalil, Asma Ambreen, and Tariq Butt. Comparison of sensitivity of quantiferon-tb gold test and tuberculin skin test in active pulmonary tuberculosis. *J Coll Physicians Surg Pak*, 23(9):633–636, 2013.
- [49] Khaled M Khalil, M Abdel-Aziz, Taymour T Nazmy, and Abdel-Badeeh M Salem. An agent-based modeling for pandemic influenza in egypt. In *Handbook on Decision Making*, pages 205–218. Springer, 2012.
- [50] Thomas C King, Mark Upfal, Andrew Gottlieb, Philip Adamo, Edward Bernacki, Chris P Kadlec, Jeffrey G Jones, Frances Humphrey-Carothers, Albert F Rielly, Pamela Drewry, et al. T-spot. tb interferon- $\gamma$  release assay performance in healthcare worker screening at nineteen us hospitals. *American journal of respiratory and critical care medicine*, 192(3):367–373, 2015.
- [51] Anton J Kleywegt, Alexander Shapiro, and Tito Homem-de Mello. The sample average approximation method for stochastic discrete optimization. *SIAM Journal on Optimization*, 12(2):479–502, 2002.
- [52] Qingxia Kong, Shan Li, Nan Liu, Chung-Piaw Teo, and Zhenzhen Yan. Appointment scheduling under schedule-dependent patient no-show behavior, 2015.
- [53] The Lancet. Covid-19: protecting health-care workers. *Lancet (London, England)*, 395(10228):922, 2020.
- [54] Gilbert Laporte and François V Louveaux. The integer l-shaped method for stochastic integer programs with complete recourse. *Operations research letters*, 13(3):133–142, 1993.
- [55] Stephen A Lauer, Kyra H Grantz, Qifang Bi, Forrest K Jones, Qulu Zheng, Hannah R Meredith, Andrew S Azman, Nicholas G Reich, and Justin Lessler. The incubation period of coronavirus disease 2019 (covid-19) from publicly reported confirmed cases: estimation and application. *Annals of internal medicine*, 172(9):577–582, 2020.
- [56] Chi Chiu Leung, Wing Cheong Yam, Pak Leung Ho, Wing Wai Yew, Chi Kuen Chan, Wing Sze Law, Shuk Nor Lee, Kwok Chiu Chang, Lai Bun Tai, and Cheuk Ming Tam. T-s pot. tb outperforms tuberculin skin test in predicting development of active tuberculosis among household contacts. *Respirology*, 20(3):496–503, 2015.

- [57] Chi Chiu Leung, Wing Cheong Yam, Wing Wai Yew, Pak Leung Ho, Cheuk Ming Tam, Wing Sze Law, Ka Fai Au, and Pui Wah Tsui. T-spot. tb outperforms tuberculin skin test in predicting tuberculosis disease. *American journal of respiratory and critical care medicine*, 182(6):834–840, 2010.
- [58] Jianzhe Luo, Vidyadhar G Kulkarni, and Serhan Ziya. Appointment scheduling under patient no-shows and service interruptions. *Manufacturing & Service Operations Management*, 14(4):670–684, 2012.
- [59] Lisa M Maillart, Julie Simmons Ivy, Scott Ransom, and Kathleen Diehl. Assessing dynamic breast cancer screening policies. *Operations Research*, 56(6):1411–1427, 2008.
- [60] Wai-Kei Mak, David P Morton, and R Kevin Wood. Monte carlo bounding techniques for determining solution quality in stochastic programs. *Operations research letters*, 24(1-2):47–56, 1999.
- [61] K McCarthy, HM McGee, and CA O’Boyle. Outpatient clinic waiting times and non-attendance as indicators of quality. *Psychology, health & medicine*, 5(3):287–293, 2000.
- [62] D Menzies, M Pai, and G Comstock. New tests for the diagnosis of latent tuberculosis infection (vol 146, pg 340, 2007). *ANNALS OF INTERNAL MEDICINE*, 146(9):688–688, 2007.
- [63] Charity G Moore, Patricia Wilson-Witherspoon, and Janice C Probst. Time and money: effects of no-shows at a family practice residency clinic. *Family Medicine-Kansas City-*, 33(7):522–527, 2001.
- [64] Guillaume A Mullie, Kevin Schwartzman, Alice Zwerling, and Dieynaba S N’Diaye. Revisiting annual screening for latent tuberculosis infection in healthcare workers: a cost-effectiveness analysis. *BMC medicine*, 15(1):104, 2017.
- [65] Martha Isabel Namén León et al. A pomdp approach to age-dependent primary screening policies for cervical cancer in colombia. Master’s thesis, Bogotá-Uniandes, 2015.
- [66] Kangqi Ng, Beng Hoong Poon, Troy Hai Kiat Puar, Jessica Li Shan Quah, Wann Jia Loh, Yu Jun Wong, Thean Yen Tan, and Jagadesan Raghuram. Covid-19 and the risk to health care workers: a case report. *Annals of internal medicine*, 2020.
- [67] Yan Ni, Ke Wang, Lindu Zhao, et al. A markov decision process model of allocating emergency medical resource among multi-priority injuries. *IJMOR*, 10(1):1–17, 2017.
- [68] Ronald Ortner. Pseudometrics for state aggregation in average reward markov decision processes. In *International Conference on Algorithmic Learning Theory*, pages 373–387. Springer, 2007.
- [69] Jonathan Patrick and Martin L Puterman. Improving resource utilization for diagnostic services through flexible inpatient scheduling: A method for improving resource utilization. *Journal of the Operational Research Society*, 58(2):235–245, 2007.
- [70] Jonathan Patrick, Martin L Puterman, and Maurice Queyranne. Dynamic multipriority patient scheduling for a diagnostic resource. *Operations research*, 56(6):1507–1525, 2008.
- [71] Vishnunarayan Girishan Prabhu, Kevin Taafe, W Hand, Caglar Caglayan, Tugce Isik, and Yongjia Song. Team based, risk adjusted staffing during a pandemic: an agent based approach. *Winter Simulation Conference*, 2020.
- [72] Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.

- [73] Xiuli Qu, Yidong Peng, Nan Kong, and Jing Shi. A two-phase approach to scheduling multi-category outpatient appointments—a case study of a women’s clinic. *Health care management science*, 16(3):197–216, 2013.
- [74] Camilla Rothe, Mirjam Schunk, Peter Sothmann, Gisela Bretzel, Guenter Froeschl, Claudia Wallrauch, Thorbjörn Zimmer, Verena Thiel, Christian Janke, Wolfgang Guggemos, et al. Transmission of 2019-ncov infection from an asymptomatic contact in germany. *New England Journal of Medicine*, 382(10):970–971, 2020.
- [75] Daniel Sanchez-Taltavull, Daniel Candinas, Edgar Roldan, and Guido Beldi. Modelling strategies to organize healthcare workforce during pandemics: application to covid-19. *medRxiv*, 2020.
- [76] Dorry L Segev, Sommer E Gentry, Daniel S Warren, Brigitte Reeb, and Robert A Montgomery. Kidney paired donation and optimizing the use of live donor organs. *Jama*, 293(15):1883–1890, 2005.
- [77] Steven M Shechter, Matthew D Bailey, Andrew J Schaefer, and Mark S Roberts. The optimal time to initiate hiv therapy under ordered health states. *Operations Research*, 56(1):20–33, 2008.
- [78] Sabine Sickinger and Rainer Kolisch. The performance of a generalized bailey–welch rule for outpatient appointment scheduling under inpatient and emergency demand. *Health care management science*, 12(4):408, 2009.
- [79] Tanu Singhal. A review of coronavirus disease-2019 (covid-19). *The Indian Journal of Pediatrics*, pages 1–6, 2020.
- [80] Barry C Smith, John F Leimkuhler, and Ross M Darrow. Yield management at american airlines. *interfaces*, 22(1):8–31, 1992.
- [81] A Spinelli and G Pellino. Covid-19 pandemic: perspectives on an unfolding crisis. *Br J Surg*, 10, 2020.
- [82] Lauren N Steimle and Brian T Denton. Markov decision processes for screening and treatment of chronic diseases. In *Markov Decision Processes in Practice*, pages 189–222. Springer, 2017.
- [83] Rhonda L Stuart and Elizabeth E Gillespie. Preparing for an influenza pandemic: healthcare workers’ opinions on working during a pandemic. *Healthcare infection*, 13(3):95–99, 2008.
- [84] Ashwin Swaminathan, Rhea Martin, Sandi Gamon, Craig Aboltins, Eugene Athan, George Braitberg, Michael G Catton, Louise Cooley, Dominic E Dwyer, Deidre Edmonds, et al. Personal protective equipment and antiviral drug use during hospitalization for suspected avian or pandemic influenza1. *Emerging infectious diseases*, 13(10):1541, 2007.
- [85] Elizabeth A Talbot, Dawn Harland, Wendy Wieland-Alter, Sherry Burrer, and Lisa V Adams. Specificity of the tuberculin skin test and the t-spot. tb assay among students in a low-tuberculosis incidence setting. *Journal of American College Health*, 60(1):94–96, 2012.
- [86] Jiafu Tang and Yu Wang. An adjustable robust optimisation method for elective and emergency surgery capacity allocation with demand uncertainty. *International Journal of Production Research*, 53(24):7317–7328, 2015.
- [87] Richard M Van Slyke and Roger Wets. L-shaped linear programs with applications to optimal control and stochastic programming. *SIAM Journal on Applied Mathematics*, 17(4):638–663, 1969.



- [88] Bram Verweij, Shabbir Ahmed, Anton J Kleywegt, George Nemhauser, and Alexander Shapiro. The sample average approximation method applied to stochastic routing problems: a computational study. *Computational Optimization and Applications*, 24(2-3):289–333, 2003.
- [89] J Wang, M Zhou, and F Liu. Reasons for healthcare workers becoming infected with novel coronavirus disease 2019 (covid-19) in china. *J Hosp Infect*, 20, 2020.
- [90] Lawrence R Weatherford and Samuel E Bodily. A taxonomy and research overview of perishable-asset revenue management: Yield management, overbooking, and pricing. *Operations research*, 40(5):831–844, 1992.
- [91] Peter Wrighton-Smith, Laurie Sneed, Frances Humphrey, Xuguang Tao, and Edward Bernacki. Screening health care workers with interferon- $\gamma$  release assay versus tuberculin skin test: impact on costs and adherence to testing (the switch study). *Journal of occupational and environmental medicine*, 54(7):806–815, 2012.
- [92] S Youakim. The occupational risk of tuberculosis in a low-prevalence population. *Occupational Medicine*, 66(6):466–470, 2016.
- [93] Gregory S Zaric and Margaret L Brandeau. Optimal investment in a portfolio of hiv prevention programs. *Medical Decision Making*, 21(5):391–408, 2001.